

504 **A Recipe for implementing and deploying our strategy**

505 We here outline more explicitly how Corollary 7 and Proposition 11 may be used to formulate a fully
506 differentiable objective by which a model may be trained.

507 First, if one wishes to make hard labels, namely $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$, it will first be necessary to use a
508 surrogate class of soft hypotheses $\mathcal{H}' \subseteq \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$ during training, before reverting to hard labels
509 for example by taking the mean label or the one with highest probability. Using soft hypotheses
510 during training is necessary to ensure that the empirical j -risks $R_S^j(Q)$ are differentiable in the model
511 parameters. Since how one chooses to do this will depend on the specific use case, we restrict our
512 attention here to the case of soft hypotheses. Specifically, we consider a class of soft hypotheses
513 $\mathcal{H} = \{h_\theta : \theta \in \mathbb{R}^N\} \subseteq \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$ parameterised by the weights $\theta \in \mathbb{R}^N$ of some neural network
514 of a given architecture with N parameters in such a way that the $R_S^j(h_\theta)$ are differentiable in θ . A
515 concrete example would be multiclass classification using a fully connected neural network with
516 output being softmax probabilities on the classes so that the $R_S^j(h_\theta)$ are differentiable.

517 Second, it is necessary to restrict the prior and posterior $P, Q \in \mathcal{M}(\mathcal{H})$ to a parameterised subset
518 of $\mathcal{M}(\mathcal{H})$ in which $\text{KL}(Q\|P)$ has a closed form which is differentiable in the parameterisation. A
519 simple choice for our case of a neural network with N parameters is $P, Q \in \{\mathcal{N}(\mathbf{w}, \text{diag}(\mathbf{s})) : \mathbf{w} \in$
520 $\mathbb{R}^N, \mathbf{s} \in \mathbb{R}_{>0}^N\}$. For prior a $P_{\mathbf{v}, \mathbf{r}} = \mathcal{N}(\mathbf{v}, \text{diag}(\mathbf{r}))$ and posterior $Q_{\mathbf{w}, \mathbf{s}} = \mathcal{N}(\mathbf{w}, \text{diag}(\mathbf{s}))$ we have
521 the closed form

$$\text{KL}(Q_{\mathbf{w}, \mathbf{s}}\|P_{\mathbf{v}, \mathbf{r}}) = \frac{1}{2} \left[\sum_{n=1}^N \left(\frac{s_n}{r_n} + \frac{(w_n - v_n)^2}{r_n} + \ln \frac{r_n}{s_n} \right) - N \right],$$

522 which is indeed differentiable in $\mathbf{v}, \mathbf{r}, \mathbf{w}$ and \mathbf{s} . While $Q_{\mathbf{w}, \mathbf{s}}$ and $P_{\mathbf{v}, \mathbf{r}}$ are technically distributions on
523 \mathbb{R}^D rather than \mathcal{H} , the KL-divergence between the distributions they induce on \mathcal{H} will be at most as
524 large as the expression above. Thus, substituting the expression above into the bounds we prove in
525 Section 3 can only increase the value of the bounds, meaning the enlarged bounds certainly still hold
526 with probability at least $1 - \delta$.

527 Third, in all but the simplest cases $R_S^j(Q_{\mathbf{w}, \mathbf{s}})$ will not have a closed form, much less one that is
528 differentiable in \mathbf{w} and \mathbf{s} . A common solution to this is to use the so-called pathwise gradient
529 estimator. In our case, this corresponds to drawing $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$, where \mathbb{I} is the $N \times N$ identity
530 matrix, and estimating

$$\nabla_{\mathbf{w}, \mathbf{s}} R_S^j(Q_{\mathbf{w}, \mathbf{s}}) = \nabla_{\mathbf{w}, \mathbf{s}} \left[\mathbb{E}_{\epsilon' \sim \mathcal{N}(\mathbf{0}, \mathbb{I})} R_S^j(h_{\mathbf{w} + \epsilon' \odot \sqrt{\mathbf{s}}}) \right] \approx \nabla_{\mathbf{w}, \mathbf{s}} R_S^j(h_{\mathbf{w} + \epsilon \odot \sqrt{\mathbf{s}}}),$$

531 where $h_{\mathbf{w}}$ denotes the function expressed by the neural network with parameters \mathbf{w} . For a proof that
532 this is an unbiased estimator, and for other methods for estimating the gradients of expectations, see
533 the survey [26].

534 Fourth, one must choose the prior. Designing priors which are optimal in some sense (*i.e.*, minimising
535 the Kullback-Leibler term in the right-hand side of generalisation bounds) has been at the core of an
536 active line of work in the PAC-Bayesian literature. For the sake of simplicity, and since it is out of the
537 scope of our contributions, we assume here that the prior is given beforehand, although we stress that
538 practitioners should pay great attention to its tuning. For our purposes, it suffices to say that if one
539 is using a data-dependent prior then it is necessary to partition the sample into $S = S_{\text{Prior}} \cup S_{\text{Bound}}$,
540 where S_{Prior} is used to train the prior and S_{Bound} is used to evaluate the bound. Since our bound holds
541 uniformly over posteriors $Q \in \mathcal{M}(\mathcal{H})$, the entire sample S is free to be used to train the posterior
542 Q . For a more in-depth discussion on the choice of prior, we refer to the following body of work:
543 Ambroladze et al. [2], Lever et al. [20, 21], Parrado-Hernández et al. [29], Dziugaite and Roy [13, 14],
544 Rivasplata et al. [32], Letarte et al. [19], Pérez-Ortiz et al. [30], Dziugaite et al. [12], Biggs and Guedj
545 [4, 6, 5].

546 Finally, given a confidence level $\delta \in (0, 1]$, one may use Algorithm 1 to obtain a posterior $Q_{\mathbf{w}, \mathbf{s}}$
547 with minimal upper bound on the total risk. Note we take the pointwise logarithm of the variances
548 \mathbf{r} and \mathbf{s} to obtain unbounded parameters on which to perform stochastic gradient descent or some
549 other minimisation algorithm. We use \oplus to denote vector concatenation. The algorithm can be
550 straightforwardly adapted to permit mini-batches by, for each epoch, sequentially repeating the steps
551 with S equal to each mini-batch.

Input:

\mathcal{X}, \mathcal{Y} /* Arbitrary input and output spaces */

$\bigcup_{j=1}^M E_j = \mathcal{Y}^2$ /* A finite partition into error types */

$\ell \in [0, \infty)^M$ /* A vector of losses, not all equal */

$S = S_{\text{Prior}} \cup S_{\text{Bound}} \in (\mathcal{X} \times \mathcal{Y})^m$ /* A partitioned i.i.d. sample */

$N \in \mathbb{N}$ /* The number of model parameters */

$P_{v,r}, v(S_{\text{Prior}}) \in \mathbb{R}^N, r(S_{\text{Prior}}) \in \mathbb{R}_{\geq 0}^N$ /* A (data-dependent) prior */

$Q_{w_0, s_0}, w_0 \in \mathbb{R}^N, s_0 \in \mathbb{R}_{\geq 0}^N$ /* An initial posterior */

$\delta \in (0, 1]$ /* A confidence level */

$\lambda > 0$ /* A learning rate */

T /* The number of epochs to train for */

Output:

$Q_{w,s}, w \in \mathbb{R}^N, s \in \mathbb{R}_{\geq 0}^N$ /* A trained posterior */

Procedure:

$\zeta_0 \leftarrow \log s_0$ /* Transform to unbounded scale parameters */

$p \leftarrow w_0 \oplus \zeta_0$ /* Collect mean and scale parameters */

for $t \leftarrow 1$ **to** T **do**

Draw $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$

$u \leftarrow R_S \left(h_{w+\epsilon \odot \sqrt{\exp(\zeta)}} \right)$

$B \leftarrow \frac{1}{m} \left[\text{KL} (Q_{w, \exp(\zeta)} \| P_{v,r}) + \ln \left(\frac{1}{\delta} \sqrt{\pi} e^{1/12m} \left(\frac{m}{2} \right)^{\frac{M-1}{2}} \sum_{z=0}^{M-1} \binom{M}{z} \frac{1}{(\pi m)^{z/2} \Gamma(\frac{M-z}{2})} \right) \right]$

$\tilde{u} \leftarrow (u_1, \dots, u_M, B)$

$G \leftarrow \mathbf{0}_{2N \times (M+1)}$ /* Initialise gradient matrix */

$F \leftarrow \mathbf{0}_{M+1}$ /* Initialise gradient vector */

for $j \leftarrow 1$ **to** $M+1$ **do**

$F_j \leftarrow \frac{\partial f_j^*}{\partial \tilde{u}_j}(\tilde{u})$ /* Gradients of total loss from Prop 11 */

for $i \leftarrow 1$ **to** $2N$ **do**

$G_{i,j} \leftarrow \frac{\partial \tilde{u}_j}{\partial p_i}(p)$ /* Gradients of empirical risks and bound */

end

end

$H \leftarrow GF$ /* Gradients of total loss w.r.t. parameters */

$p \leftarrow p - \lambda H$ /* Gradient step */

end

$w = (p_1, \dots, p_N)$

$s = (\exp(p_{N+1}), \dots, \exp(p_{2N}))$

return w, s

Algorithm 1: Calculating a posterior with minimal bound on the total risk.

552 B Proofs

553 B.1 Proof of Lemma 5

554 Let $\mathbf{E}_M := \{e_1, \dots, e_M\}$, namely the set of M -dimensional basis vectors. We will denote a typical
 555 element of \mathbf{E}_M^m by $\boldsymbol{\eta}^{(m)} = (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_m)$. For any $\mathbf{x}^{(m)} = (\mathbf{x}_1, \dots, \mathbf{x}_m) \in \Delta_M^m$, a straightforward
 556 induction on m yields

$$\sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) = 1. \quad (8)$$

557 To see this, for $m = 1$ we have $\mathbf{E}_M^1 = \{(e_1,), \dots, (e_M,)\}$, where we have been pedantic in using
 558 1-tuples to maintain consistency with larger values of m . Thus, for any $\mathbf{x}^{(1)} = (\mathbf{x}_1,) \in \Delta_M^1$, the left

559 hand side of equation (8) can be written as

$$\sum_{j=1}^M \mathbf{x}_1 \cdot \mathbf{e}_j = \sum_{j=1}^M (\mathbf{x}_1)_j = 1.$$

560 Now suppose that equation (8) holds for any $\mathbf{x}^{(m)} \in \Delta_M^m$ and let $\mathbf{x}^{(m+1)} = (\mathbf{x}_1, \dots, \mathbf{x}_{m+1}) \in$
 561 Δ_M^{m+1} . Then the left hand side of equation (8) can be written as

$$\begin{aligned} \sum_{\boldsymbol{\eta}^{(m+1)} \in \mathbf{E}_M^{m+1}} \left(\prod_{i=1}^{m+1} \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) &= \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \sum_{j=1}^M \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) (\mathbf{x}_{m+1} \cdot \mathbf{e}_j) \\ &= \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) \sum_{j=1}^M (\mathbf{x}_{m+1} \cdot \mathbf{e}_j) = 1. \end{aligned}$$

562 We now show that any $\mathbf{x}^{(m)} = (\mathbf{x}_1, \dots, \mathbf{x}_m) \in \Delta_M^m$ can be written as a convex combination of the
 563 elements of \mathbf{E}_M^m in the following way

$$\mathbf{x}^{(m)} = \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) \boldsymbol{\eta}^{(m)}. \quad (9)$$

564 We have already shown that the weights sum to one, and they are clearly elements of $[0, 1]$, so the
 565 right hand side of equation (9) is indeed a convex combination of the elements of \mathbf{E}_M^m . We now show
 566 that equation (9) holds, again by induction.

567 For $m = 1$ and any $\mathbf{x}^{(1)} = (\mathbf{x}_1,) \in \Delta_M^1$, the right hand side of equation (9) can be written as

$$\sum_{j=1}^M (\mathbf{x}_1 \cdot \mathbf{e}_j) (\mathbf{e}_j,) = (\mathbf{x}_1,) = \mathbf{x}.$$

568 For the inductive hypothesis, suppose equation (9) holds for some arbitrary $m \geq 1$, and denote
 569 elements of \mathbf{E}_M^{m+1} by $\boldsymbol{\eta}^{(m)} \oplus (\mathbf{e},)$ for some $\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m$ and $\mathbf{e} \in \mathbf{E}_M$, where \oplus denotes vector
 570 concatenation. Then for any $\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} \oplus (\mathbf{x}_{m+1},) = (\mathbf{x}_1, \dots, \mathbf{x}_{m+1}) \in \Delta_M^{m+1}$, the right
 571 hand side of equation (9) can be written as

$$\begin{aligned} \sum_{\boldsymbol{\eta}^{(m+1)} \in \mathbf{E}_M^{m+1}} \left(\prod_{i=1}^{m+1} \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) \boldsymbol{\eta}^{(m+1)} &= \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \sum_{j=1}^M \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) (\mathbf{x}_{m+1} \cdot \mathbf{e}_j) \boldsymbol{\eta}^{(m)} \oplus (\mathbf{e}_j,) \\ &= \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \sum_{j=1}^M \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) (\mathbf{x}_{m+1} \cdot \mathbf{e}_j) \boldsymbol{\eta}^{(m)} \\ &\quad \oplus \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \sum_{j=1}^M \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) (\mathbf{x}_{m+1} \cdot \mathbf{e}_j) (\mathbf{e}_j,) \\ &= \sum_{j=1}^M (\mathbf{x}_{m+1} \cdot \mathbf{e}_j) \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) \boldsymbol{\eta}^{(m)} \\ &\quad \oplus \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i \right) \sum_{j=1}^M (\mathbf{x}_{m+1} \cdot \mathbf{e}_j) (\mathbf{e}_j,) \\ &= 1 \cdot \mathbf{x}^{(m)} \oplus 1 \cdot (\mathbf{x}_{m+1},) = \mathbf{x}^{(m+1)}, \end{aligned}$$

572 where in the penultimate equality we have used the inductive hypothesis and (twice) the result of the
 573 previous induction.

574 We can now prove the statement of the Lemma. Applying Jensen's inequality to equation (9) with the
 575 convex function f , we have that

$$\begin{aligned} f(\mathbf{x}_1, \dots, \mathbf{x}_m) &= f\left(\sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i\right) \boldsymbol{\eta}^{(m)}\right) \\ &\leq \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \left(\prod_{i=1}^m \mathbf{x}_i \cdot \boldsymbol{\eta}_i\right) f(\boldsymbol{\eta}^{(m)}). \end{aligned}$$

576 Let $\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}_1]$ denote the mean of the i.i.d. random vectors X_i . Then the above inequality implies

$$\begin{aligned} \mathbb{E}[f(\mathbf{X}_1, \dots, \mathbf{X}_m)] &\leq \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \left(\prod_{i=1}^m \boldsymbol{\mu} \cdot \boldsymbol{\eta}_i\right) f(\boldsymbol{\eta}^{(m)}) \\ &= \sum_{\boldsymbol{\eta}^{(m)} \in \mathbf{E}_M^m} \left(\prod_{i=1}^m \mathbb{P}(\mathbf{X}'_i = \boldsymbol{\eta}_i)\right) f(\boldsymbol{\eta}^{(m)}) \\ &= \mathbb{E}[f(\mathbf{X}'_1, \dots, \mathbf{X}'_m)]. \end{aligned}$$

577 B.2 Proof of Lemma 8

578 The proof of Lemma 8 itself requires two technical helping lemmas which we now state and prove.

579 **Lemma 12.** For any integers $n \geq 2$ and $p \geq -1$,

$$\sum_{k=1}^{n-1} \frac{(n-k)^{p/2}}{\sqrt{k}} \leq n^{\frac{p+1}{2}} \int_0^1 \frac{(1-x)^{p/2}}{\sqrt{x}} dx.$$

580 *Proof.* The case of $p = -1$, namely

$$\sum_{k=1}^{n-1} \frac{1}{\sqrt{k(n-k)}} \leq \int_0^1 \frac{1}{\sqrt{x(1-x)}} dx,$$

581 has already been demonstrated in [22]. For $p > -1$, let

$$f_p(x) := \frac{(1-x)^{p/2}}{\sqrt{x}}.$$

582 We will show that each $f_p(\cdot)$ is monotonically decreasing on $(0, 1)$. Indeed,

$$\frac{df_p}{dx}(x) = -\frac{(1-x)^{\frac{p}{2}-1}(px+1-x)}{2x^{3/2}} \leq -\frac{(1-x)^{p/2}}{2x^{3/2}} < 0,$$

583 where for the inequalities we have used the fact that $p > -1$ and $x \in (0, 1)$. We therefore see that

$$\begin{aligned} \sum_{k=1}^{n-1} \frac{(n-k)^{p/2}}{\sqrt{k}} &= \sum_{k=1}^{n-1} \frac{n^{p/2} (1 - \frac{k}{n})^{p/2}}{\sqrt{n} \sqrt{\frac{k}{n}}} \\ &= n^{\frac{p+1}{2}} \sum_{k=1}^{n-1} \frac{1}{n} \frac{(1 - \frac{k}{n})^{p/2}}{\sqrt{\frac{k}{n}}} \\ &= n^{\frac{p+1}{2}} \sum_{k=1}^{n-1} \frac{1}{n} f_p\left(\frac{k}{n}\right) \\ &\leq n^{\frac{p+1}{2}} \sum_{k=1}^{n-1} \int_{\frac{k-1}{n}}^{\frac{k}{n}} f_p(x) dx \end{aligned}$$

$$\begin{aligned}
&= n^{\frac{p+1}{2}} \int_0^{1-\frac{1}{n}} f_p(x) dx \\
&\leq n^{\frac{p+1}{2}} \int_0^1 f_p(x) dx.
\end{aligned}$$

584

□

585 Intuitively, the proof of the above lemma works by bounding the integral below by a Riemann sum.
586 In the following lemma we actually calculate this integral, yielding a more explicit bound on the sum
587 in Lemma 12. We found it is easier to calculate a slightly more general integral, where the 1 in the
588 limit and the integrand is replaced by a positive constant a .

589 **Lemma 13.** For any real number $a > 0$ and integer $n \geq -1$,

$$\int_0^a \frac{(a-x)^{n/2}}{\sqrt{x}} dx = \sqrt{\pi} \frac{\Gamma(\frac{n+2}{2})}{\Gamma(\frac{n+3}{2})} a^{\frac{n+1}{2}}.$$

590 *Proof.* Define

$$I_n(a) := \int_0^a \frac{(a-x)^{n/2}}{\sqrt{x}} dx \quad \text{and} \quad f_n(a) := \sqrt{\pi} \frac{\Gamma(\frac{n+2}{2})}{\Gamma(\frac{n+3}{2})} a^{\frac{n+1}{2}}.$$

591 We proceed by induction, increasing n by 2 each time. This means we need two base cases. First, for
592 $n = -1$, we have

$$I_{-1}(a) = \int_0^a \frac{1}{\sqrt{x(a-x)}} dx = \left[2 \arcsin \sqrt{\frac{x}{a}} \right]_0^a = \pi = f_{-1}(a),$$

593 since $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ and $\Gamma(1) = 1$. Second, for $n = 0$,

$$I_0(a) = \int_0^a \frac{1}{\sqrt{x}} dx = [2\sqrt{x}]_0^a = 2\sqrt{a} = f_0(a),$$

594 since $\Gamma(\frac{3}{2}) = \frac{\sqrt{\pi}}{2}$. Now, by the Leibniz integral rule, we have

$$\frac{d}{da} I_{n+2}(a) = \int_0^a \frac{\partial}{\partial a} \frac{(a-x)^{\frac{n+2}{2}}}{\sqrt{x}} dx = \frac{n+2}{2} \int_0^a \frac{(a-x)^{\frac{n}{2}}}{\sqrt{x}} dx = \frac{n+2}{2} I_n(a).$$

595 Thus

$$I_{n+2}(a) = \frac{n+2}{2} \left[\int_0^a I_n(t) dt + I_n(0) \right] = \frac{n+2}{2} \int_0^a I_n(t) dt,$$

596 since $I_n(0) = 0$.

597 Now, for the inductive step, suppose $I_n(a) = f_n(a)$ for some $n \geq -1$. Then, using the previous
598 calculation, we have

$$\begin{aligned}
I_{n+2}(a) &= \frac{n+2}{2} \int_0^a f_n(t) dt \\
&= \frac{n+2}{2} \int_0^a \sqrt{\pi} \frac{\Gamma(\frac{n+2}{2})}{\Gamma(\frac{n+3}{2})} t^{\frac{n+1}{2}} dt \\
&= \sqrt{\pi} \frac{\frac{n+2}{2} \Gamma(\frac{n+2}{2})}{\frac{n+3}{2} \Gamma(\frac{n+3}{2})} a^{\frac{n+3}{2}} \\
&= \sqrt{\pi} \frac{\Gamma(\frac{n+2}{2} + 1)}{\Gamma(\frac{n+3}{2} + 1)} a^{\frac{n+3}{2}} \\
&= \sqrt{\pi} \frac{\Gamma\left(\frac{(n+2)+2}{2}\right)}{\Gamma\left(\frac{(n+2)+3}{2}\right)} a^{\frac{(n+2)+1}{2}} \\
&= f_{n+2}(a).
\end{aligned}$$

599 This completes the proof.

□

600 We are now ready to prove Lemma 8 which, for ease of reference, we restate here. For integers
 601 $M \geq 1$ and $m \geq M$,

$$\sum_{\mathbf{k} \in S_{m,M}^{>0}} \frac{1}{\prod_{j=1}^M \sqrt{k_j}} \leq \frac{\pi^{\frac{M}{2}} m^{\frac{M-2}{2}}}{\Gamma(\frac{M}{2})}.$$

602 *Proof.* (of Lemma 8) We proceed by induction on M . For $M = 1$, the set $S_{m,M}$ contains a single
 603 element, namely the one-dimensional vector $\mathbf{k} = (k_1, \dots) = (m, \dots)$. In this case, the left hand side is
 604 $1/\sqrt{m}$ while the right hand side is $\sqrt{\pi}/(\sqrt{m}\Gamma(1/2)) = 1/\sqrt{m}$, since $\Gamma(1/2) = \sqrt{\pi}$.

605 Now, as the inductive hypothesis, assume the inequality of Lemma 8 holds for some fixed $M \geq 1$
 606 and all $m \geq M$. Then for all $m \geq M + 1$, we have

$$\begin{aligned} \sum_{\mathbf{k} \in S_{m,M+1}^{>0}} \frac{1}{\prod_{j=1}^{M+1} \sqrt{k_j}} &= \sum_{k_1=1}^{m-M} \frac{1}{\sqrt{k_1}} \sum_{\mathbf{k}' \in S_{m-k_1,M}^{>0}} \frac{1}{\prod_{j=1}^M \sqrt{k'_j}} \\ &\leq \sum_{k_1=1}^{m-M} \frac{1}{\sqrt{k_1}} \frac{\pi^{\frac{M}{2}} (m-k_1)^{\frac{M-2}{2}}}{\Gamma(\frac{M}{2})} \quad (\text{by the inductive hypothesis}) \\ &= \frac{\pi^{\frac{M}{2}}}{\Gamma(\frac{M}{2})} \sum_{k_1=1}^{m-M} \frac{(m-k_1)^{\frac{M-2}{2}}}{\sqrt{k_1}} \\ &\leq \frac{\pi^{\frac{M}{2}}}{\Gamma(\frac{M}{2})} \sum_{k_1=1}^{m-1} \frac{(m-k_1)^{\frac{M-2}{2}}}{\sqrt{k_1}} \quad (\text{enlarging the sum domain}) \\ &\leq \frac{\pi^{\frac{M}{2}}}{\Gamma(\frac{M}{2})} m^{\frac{M-1}{2}} \int_0^1 \frac{(1-x)^{\frac{M-2}{2}}}{\sqrt{x}} dx \quad (\text{by Lemma 12}) \\ &= \frac{\pi^{\frac{M}{2}}}{\Gamma(\frac{M}{2})} m^{\frac{M-1}{2}} \sqrt{\pi} \frac{\Gamma(\frac{M}{2})}{\Gamma(\frac{M+1}{2})} \quad (\text{by Lemma 13}) \\ &= \frac{\pi^{\frac{M+1}{2}} m^{\frac{M-1}{2}}}{\Gamma(\frac{M+1}{2})}, \end{aligned}$$

607 as required. □

608 B.3 Proof of Proposition 9

609 *Proof.* The case where $q_j = 1$ or $p_j = 1$ can be dealt with trivially by splitting into the three
 610 following subcases

- 611 • $q_j = p_j = 1 \implies \text{kl}(q_j \| p_j) = \text{kl}(\mathbf{q} \| \mathbf{p}) = 0$
- 612 • $q_j = 1, p_j \neq 1 \implies \text{kl}(q_j \| p_j) = \text{kl}(\mathbf{q} \| \mathbf{p}) = -\log p_j$
- 613 • $q_j \neq 1, p_j = 1 \implies \text{kl}(q_j \| p_j) = \text{kl}(\mathbf{q} \| \mathbf{p}) = \infty$.

614 For $q_j \neq 1$ and $p_j \neq 1$ define the distributions $\tilde{\mathbf{q}}, \tilde{\mathbf{p}} \in \Delta_M$ by $\tilde{q}_j = \tilde{p}_j = 0$ and

$$\tilde{q}_i = \frac{q_i}{1-q_j} \quad \text{and} \quad \tilde{p}_i = \frac{p_i}{1-p_j}$$

615 for $i \neq j$. Then

$$\begin{aligned} \sum_{i \neq j} q_i \log \frac{q_i}{p_i} &= \sum_{i \neq j} (1-q_j) \tilde{q}_i \log \frac{(1-q_j) \tilde{q}_i}{(1-p_j) \tilde{p}_i} \\ &= (1-q_j) \sum_{i \neq j} \tilde{q}_i \log \frac{\tilde{q}_i}{\tilde{p}_i} + \tilde{q}_i \log \frac{1-q_j}{1-p_j} \end{aligned}$$

$$\begin{aligned}
&= (1 - q_j) \text{kl}(\tilde{\mathbf{q}} \parallel \tilde{\mathbf{p}}) + (1 - q_j) \log \frac{1 - q_j}{1 - p_j} \\
&\geq (1 - q_j) \log \frac{1 - q_j}{1 - p_j}.
\end{aligned}$$

616 The final inequality holds since $\text{kl}(\tilde{\mathbf{q}} \parallel \tilde{\mathbf{p}}) \geq 0$. Further, note that we have equality if and only if $\tilde{\mathbf{q}} = \tilde{\mathbf{p}}$,
617 which, by their definitions, translates to

$$p_i = \frac{1 - p_j}{1 - q_j} q_i$$

618 for all $i \neq j$. If we now add $q_j \log \frac{q_j}{p_j}$ to both sides, we obtain

$$\text{kl}(\mathbf{q} \parallel \mathbf{p}) \geq (1 - q_j) \log \frac{1 - q_j}{1 - p_j} + q_j \log \frac{q_j}{p_j} = \text{kl}(q_j \parallel p_j),$$

619 with the same condition for equality. \square

620 The following proposition makes more precise the argument found at the beginning of Section 4
621 for how Proposition 9 can be used to derive the tightest possible lower and upper bounds on each
622 $R_D^j(Q)$.

623 **Proposition 14.** *Suppose that $\mathbf{q}, \mathbf{p} \in \Delta_M$ are such that $\text{kl}(\mathbf{q} \parallel \mathbf{p}) \leq B$, where \mathbf{q} is known and \mathbf{p} is
624 unknown. Then, in the absence of any further information, the tightest bound that can be obtained on
625 each p_j is*

$$p_j \leq \text{kl}^{-1}(q_j, B).$$

626 *Proof.* Suppose $p_j > \text{kl}^{-1}(q_j, B)$. Then, by definition of kl^{-1} , we have that $\text{kl}(q_j \parallel p_j) > B$.
627 By Proposition 9, this would then imply $\text{kl}(\mathbf{q} \parallel \mathbf{p}) > B$, contradicting our assumption. Therefore
628 $p_j \leq \text{kl}^{-1}(q_j, B)$. Now, with the information we have, we cannot rule out that

$$p_i = \frac{1 - p_j}{1 - q_j} q_i$$

629 for all $i \neq j$ and thus, by Proposition 9, that $\text{kl}(q_j \parallel p_j) = \text{kl}(\mathbf{q} \parallel \mathbf{p})$. Further, we cannot rule out that
630 $\text{kl}(\mathbf{q} \parallel \mathbf{p}) = B$. Thus, it is possible that $\text{kl}(q_j \parallel p_j) = B$, in which case $p_j = \text{kl}^{-1}(q_j, B)$. We therefore
631 see that $\text{kl}^{-1}(q_j, B)$ is the tightest possible upper bound on p_j , for each $j \in [M]$. \square

632 B.4 Proof of Proposition 11

633 Before proving the proposition, we first argue that $\text{kl}_\ell^{-1}(\mathbf{u} \mid c)$ given by Definition 10 is well-defined.
634 First, note that $A_{\mathbf{u}} := \{\mathbf{v} \in \Delta_M : \text{kl}(\mathbf{u} \parallel \mathbf{v}) \leq c\}$ is compact (boundedness is clear and it is closed
635 because it is the preimage of the closed set $[0, c]$ under the continuous map $\mathbf{v} \mapsto \text{kl}(\mathbf{u} \parallel \mathbf{v})$) and so the
636 continuous function f_ℓ achieves its supremum on $A_{\mathbf{u}}$. Further, note that $A_{\mathbf{u}}$ is a convex subset of
637 Δ_M (because the map $\mathbf{v} \mapsto \text{kl}(\mathbf{u} \parallel \mathbf{v})$ is convex) and f_ℓ is linear, so the supremum of f_ℓ over $A_{\mathbf{u}}$
638 is achieved and is located on the boundary of $A_{\mathbf{u}}$. This means we can replace the inequality constraint
639 $\text{kl}(\mathbf{u} \parallel \mathbf{v}) \leq c$ in Definition 10 with the equality constraint $\text{kl}(\mathbf{u} \parallel \mathbf{v}) = c$. Finally, if $\mathbf{u} \in \Delta_M^{>0}$ then
640 $A_{\mathbf{u}}$ is a *strictly* convex subset of Δ_M (because the map $\mathbf{v} \mapsto \text{kl}(\mathbf{u} \parallel \mathbf{v})$ is then *strictly* convex) and so
641 the supremum of f_ℓ occurs at a *unique* point on the boundary of $A_{\mathbf{u}}$. In other words, if $\mathbf{u} \in \Delta_M^{>0}$
642 then $\text{kl}_\ell^{-1}(\mathbf{u} \mid c)$ is defined *uniquely*.

643 *Proof.* (of Proposition 11) We start by deriving the implicit expression for $\mathbf{v}^*(\tilde{\mathbf{u}}) = \text{kl}_\ell^{-1}(\mathbf{u} \mid c)$ given
644 in the proposition by solving a transformed version of the optimisation problem given by Definition
645 10 using the method of Lagrange multipliers. We obtain two solutions to the Lagrangian equations,
646 which must correspond to the maximum and minimum total risk over the set $A_{\mathbf{u}} := \{\mathbf{v} \in \Delta_M : \text{kl}(\mathbf{u} \parallel \mathbf{v}) \leq c\}$
647 because, as argued in the main text (see the discussion after Definition 10), $A_{\mathbf{u}}$ is
648 compact and so the linear total risk $f_\ell(\mathbf{v})$ attains its maximum and minimum on $A_{\mathbf{u}}$.

649 By definition of $\mathbf{v}^*(\tilde{\mathbf{u}}) = \text{kl}_\ell^{-1}(\mathbf{u} \mid c)$, we know that $\text{kl}(\mathbf{v}^*(\tilde{\mathbf{u}}) \parallel \mathbf{u}) \leq c$. Since, by assumption,
650 $u_j > 0$ for all j , we see that $\mathbf{v}^*(\tilde{\mathbf{u}})_j > 0$ for all j , otherwise we would have $\text{kl}(\mathbf{v}^*(\tilde{\mathbf{u}}) \parallel \mathbf{u}) = \infty$, a

651 contradiction. Thus $\mathbf{v}^*(\tilde{\mathbf{u}}) \in \Delta_M^{\geq 0}$ and we are permitted to instead optimise over the unbounded
652 variable $\mathbf{t} \in \mathbb{R}^M$, where $t_j := \ln v_j$. With this transformation, the constraint $\mathbf{v} \in \Delta_M$ can be
653 replaced simply by $\sum_j e^{t_j} = 1$ and the optimisation problem becomes

$$\begin{aligned} \text{Maximise: } F(\mathbf{t}) &:= \sum_{j=1}^M \ell_j e^{t_j} \\ \text{Subject to: } g(\mathbf{t}; \mathbf{u}, c) &:= \text{kl}(\mathbf{u} \| e^{\mathbf{t}}) - c = 0, \\ h(\mathbf{t}) &:= \sum_{j=1}^M e^{t_j} - 1 = 0, \end{aligned}$$

654 where $e^{\mathbf{t}} \in \mathbb{R}^M$ is defined by $(e^{\mathbf{t}})_j := e^{t_j}$. Note that $F(\mathbf{t}) = f_{\ell}(e^{\mathbf{t}})$. Following the terminology
655 of mathematical economics, we call the t_j the *optimisation variables*, and the \tilde{u}_j (namely the u_j
656 and c) the *choice variables*. The vector ℓ is considered fixed—we neither want to optimise over
657 it nor differentiate with respect to it—which is why we occasionally suppress it from the notation
658 henceforth.

659 For each $\tilde{\mathbf{u}}$, let $\mathbf{v}^*(\tilde{\mathbf{u}})$ and $\mathbf{t}^*(\tilde{\mathbf{u}})$ be the solutions to the original and transformed optimisation
660 problems respectively. Since the map $\mathbf{v} = e^{\mathbf{t}}$ is one-to-one, it is clear that since $\mathbf{v}^*(\tilde{\mathbf{u}})$ exists uniquely,
661 so does $\mathbf{t}^*(\tilde{\mathbf{u}})$, and that they are related by $\mathbf{v}^*(\tilde{\mathbf{u}}) = e^{\mathbf{t}^*(\tilde{\mathbf{u}})}$. We therefore have the identity

$$f_{\ell}(\mathbf{v}^*(\tilde{\mathbf{u}})) \equiv F(\mathbf{t}^*(\tilde{\mathbf{u}})).$$

662 Recalling that $f_{\ell}^*(\tilde{\mathbf{u}}) := f_{\ell}(\mathbf{v}^*(\tilde{\mathbf{u}}))$, we see that

$$\nabla_{\tilde{\mathbf{u}}} f_{\ell}^*(\tilde{\mathbf{u}}) \equiv \nabla_{\tilde{\mathbf{u}}} F(\mathbf{t}^*(\tilde{\mathbf{u}})). \quad (10)$$

663 the derivatives of $f_{\ell}(\text{kl}_{\ell}^{-1}(\mathbf{u}|c))$ with respect to \mathbf{u} and c are given by $\nabla_{\tilde{\mathbf{u}}} F(\mathbf{t}^*(\tilde{\mathbf{u}}))$.

664 Using the method of Lagrange multipliers, there exist real numbers $\lambda^* = \lambda^*(\tilde{\mathbf{u}})$ and $\mu^* = \mu^*(\tilde{\mathbf{u}})$
665 such that $(\mathbf{t}^*, \lambda^*, \mu^*)$ is a stationary point (with respect to \mathbf{t} , λ and μ) of the Lagrangian function

$$\mathcal{L}(\mathbf{t}, \lambda, \mu; \tilde{\mathbf{u}}) := F(\mathbf{t}) + \lambda g(\mathbf{t}; \tilde{\mathbf{u}}) + \mu h(\mathbf{t}).$$

666 Let $F_{\mathbf{t}}(\cdot)$ and $h_{\mathbf{t}}(\cdot)$ denote the gradient vectors of F and h respectively, and let $g_{\mathbf{t}}(\cdot; \tilde{\mathbf{u}})$ and $g_{\tilde{\mathbf{u}}}(\mathbf{t}; \cdot)$
667 denote the gradient vectors of g with respect to \mathbf{t} only and $\tilde{\mathbf{u}}$ only, respectively. Simple calculation
668 yields

$$\begin{aligned} g_{\mathbf{t}}(\mathbf{t}; \tilde{\mathbf{u}}) &= \left(\frac{\partial g}{\partial t_1}(\mathbf{t}; \tilde{\mathbf{u}}), \dots, \frac{\partial g}{\partial t_M}(\mathbf{t}; \tilde{\mathbf{u}}) \right) = -\mathbf{u} \quad \text{and} \\ g_{\tilde{\mathbf{u}}}(\mathbf{t}; \tilde{\mathbf{u}}) &= \left(\frac{\partial g}{\partial \tilde{u}_1}(\mathbf{t}; \tilde{\mathbf{u}}), \dots, \frac{\partial g}{\partial \tilde{u}_{M+1}}(\mathbf{t}; \tilde{\mathbf{u}}) \right) = \left(1 - t_1 + \log u_1, \dots, 1 - t_M + \log u_M, -1 \right). \end{aligned} \quad (11)$$

669 Then, taking the partial derivatives of \mathcal{L} with respect to λ, μ and the t_j , we have that $(\mathbf{t}, \lambda, \mu) =$
670 $(\mathbf{t}^*(\tilde{\mathbf{u}}), \lambda^*(\tilde{\mathbf{u}}), \mu^*(\tilde{\mathbf{u}}))$ solves the simultaneous equations

$$F_{\mathbf{t}}(\mathbf{t}) + \lambda g_{\mathbf{t}}(\mathbf{t}; \tilde{\mathbf{u}}) + \mu h_{\mathbf{t}}(\mathbf{t}) = \mathbf{0}, \quad (12)$$

671

$$g(\mathbf{t}; \tilde{\mathbf{u}}) = 0, \quad \text{and}$$

672

$$h(\mathbf{t}) = 0,$$

673 where the last two equations recover the constraints. Substituting the gradients $F_{\mathbf{t}}, g_{\mathbf{t}}$ and $h_{\mathbf{t}}$, the first
674 equation reduces to

$$\ell \odot e^{\mathbf{t}} - \lambda \mathbf{u} + \mu e^{\mathbf{t}} = \mathbf{0},$$

675 which implies that for all $j \in [M]$

$$e^{t_j} = \frac{\lambda u_j}{\mu + \ell_j}. \quad (13)$$

676 Substituting this into the constraints $g = h = 0$ yields the following simultaneous equations in λ and
677 μ

$$c = \text{kl}(\mathbf{u} \| e^{\mathbf{t}}) = \sum_{j=1}^M u_j \log \frac{u_j}{e^{t_j}} = \sum_{j=1}^M u_j \log \frac{\mu + \ell_j}{\lambda} \quad \text{and} \quad \lambda \sum_{j=1}^M \frac{u_j}{\mu + \ell_j} = 1.$$

678 Substituting the second into the first and rearranging the second, this is equivalent to solving

$$c = \sum_{j=1}^M u_j \log \left((\mu + \ell_j) \sum_{k=1}^M \frac{u_k}{\mu + \ell_k} \right) \quad \text{and} \quad \lambda = \left(\sum_{j=1}^M \frac{u_j}{\mu + \ell_j} \right)^{-1}. \quad (14)$$

679 It has already been established in the discussion after Definition 10 that $f_{\ell}(\mathbf{v})$ attains its maximum
680 on the set $A_{\mathbf{u}} := \{\mathbf{v} \in \Delta_M : \text{kl}(\mathbf{u} \parallel \mathbf{v}) \leq c\}$. Therefore $F(\mathbf{t})$ also attains its maximum on \mathbb{R}^M and
681 one of the solutions to these simultaneous equations corresponds to this maximum. We first show
682 that there is a single solution to the first equation in the set $(-\infty, -\max_j \ell_j)$, referred to as $\mu^*(\tilde{\mathbf{u}})$ in
683 the proposition. Second, we show that any other solution corresponds to a smaller total risk, so that
684 $\mu^*(\tilde{\mathbf{u}})$ corresponds to the maximum total risk and yields $\mathbf{v}^*(\tilde{\mathbf{u}}) = \text{kl}_{\ell}^{-1}(\mathbf{u} \parallel c)$ when $\mu^*(\tilde{\mathbf{u}})$ and the
685 associated $\lambda^*(\tilde{\mathbf{u}})$ are substituted into Equation 13.

686 For the first step, note that since the e^{ℓ_j} are probabilities, we see from Equation 13 that either
687 $\mu + \ell_j > 0$ for all j (in the case that $\lambda > 0$), or $\mu + \ell_j < 0$ for all j (in the case that $\lambda < 0$).
688 Thus any solutions μ to the first equation must be in $(-\infty, -\max_j \ell_j)$ or $(-\min_j \ell_j, \infty)$. If
689 $\mu \in (-\infty, -\max_j \ell_j)$ then the first equation can be written as $c = \phi_{\ell}(\mu)$, with ϕ_{ℓ} as defined in the
690 statement of the proposition. We now show that ϕ_{ℓ} is strictly increasing in μ , and that $\phi_{\ell}(\mu) \rightarrow 0$ as
691 $\mu \rightarrow -\infty$ and $\phi_{\ell}(\mu) \rightarrow \infty$ as $\mu \rightarrow -\max_j \ell_j$, so that $c = \phi_{\ell}(\mu)$ does indeed have a single solution
692 in the set $(-\infty, -\max_j \ell_j)$. Straightforward differentiation and algebra shows that

$$\begin{aligned} \phi'_{\ell}(\mu) &= \sum_{j=1}^M \frac{u_j}{(\mu + \ell_j) \sum_{k=1}^M \frac{u_k}{\mu + \ell_k}} \left(\sum_{k'=1}^M \frac{u_{k'}}{\mu + \ell_{k'}} - (\mu + \ell_j) \sum_{k'=1}^M \frac{u_{k'}}{(\mu + \ell_{k'})^2} \right) \\ &= \frac{\left(\sum_{j=1}^M \frac{u_j}{\mu + \ell_j} \right)^2 - \sum_{j=1}^M \frac{u_j}{(\mu + \ell_j)^2}}{\sum_{k=1}^M \frac{u_k}{\mu + \ell_k}}. \end{aligned}$$

693 Jensen's inequality demonstrates that the numerator is strictly negative, where strictness is due to
694 the assumption that the ℓ_j are not all equal. Further, since the denominator is strictly negative (since
695 we are dealing with the case where $\mu \in (-\infty, -\max_j \ell_j)$), we see that ϕ_{ℓ} is strictly increasing for
696 $\mu \in (-\infty, -\max_j \ell_j)$.² Turning to the limits, we first show that $\phi_{\ell}(\mu) \rightarrow \infty$ as $\mu \rightarrow -\max_j \ell_j$.

697 We now determine the left hand limit. Define $J = \{j \in [M] : \ell_j = \max_k \ell_k\}$, noting that
698 this is a strict subset of $[M]$ since by assumption the ℓ_j are not all equal. We then have that for
699 $\mu \in (-\infty, \max_j \ell_j)$

$$\begin{aligned} e^{\phi_{\ell}(\mu)} &= \left(-\sum_{j=1}^M \frac{u_j}{\mu + \ell_j} \right) \left(\prod_{k=1}^M (-\mu + \ell_k)^{u_k} \right) \\ &= \left(-\sum_{j \in J} \frac{u_j}{\mu + \ell_j} - \sum_{j' \notin J} \frac{u_{j'}}{\mu + \ell_{j'}} \right) \prod_{k \in J} (-\mu + \ell_k)^{u_k} \prod_{k' \notin J} (-\mu + \ell_{k'})^{u_{k'}} \\ &\geq \left(-\sum_{j \in J} \frac{u_j}{\mu + \ell_j} \right) \prod_{k \in J} (-\mu + \ell_k)^{u_k} \prod_{k' \notin J} (-\mu + \ell_{k'})^{u_{k'}} \\ &= \frac{\left(\sum_{j \in J} u_j \right) \left(\prod_{k' \notin J} (-\mu + \ell_{k'})^{u_{k'}} \right)}{\left(-(\mu + \max_j \ell_j) \right)^{1 - \sum_{k \in J} u_k}}. \end{aligned}$$

700 The first term in the numerator is a positive constant, independent of μ . The second term in the
701 numerator tends to a finite positive limit as $\mu \uparrow -\max_j \ell_j$. Since $[M] \setminus J$ is non-empty, the power
702 in the denominator is positive and the term in the outer brackets is positive and tends to zero as
703 $\mu \uparrow -\max_j \ell_j$. Thus $e^{\phi_{\ell}(\mu)} \rightarrow \infty$ as $\mu \uparrow -\max_j \ell_j$ and, by the continuity of the logarithm, $\phi_{\ell}(\mu)$
704 as $\mu \uparrow -\max_j \ell_j$.

²Incidentally, this argument also shows that there is at most one solution to the first equation in (14) in the range $(-\min_j \ell_j, \infty)$. There indeed exists a unique solution, which corresponds to the minimum total risk, but we do not prove this.

705 We now determine $\lim_{\mu \rightarrow -\infty} \phi_{\ell}(\mu)$ by sandwiching $\phi(\mu)$ between two functions that both tend to
 706 zero as $\mu \rightarrow -\infty$. First, since $\ell_j \geq 0$ for all j , for $\mu \in (-\infty, -\max_j \ell_j)$ we have

$$\log \left(-\sum_{j=1}^M \frac{u_j}{\mu + \ell_j} \right) \geq \log \left(-\sum_{j=1}^M \frac{u_j}{\mu} \right) = -\log(-\mu) = -\sum_{j=1}^M u_j \log(-\mu),$$

707 and so

$$\phi_{\ell}(\mu) \geq -\sum_{j=1}^M u_j \log(-\mu) + \sum_{j=1}^M u_j \log(-(\mu + \ell_j)) = \sum_{j=1}^M u_j \log \left(1 + \frac{\ell_j}{\mu} \right) \rightarrow 0 \quad \text{as } \mu \rightarrow -\infty.$$

708 Similarly,

$$\sum_{j=1}^M u_j \log(-(\mu + \ell_j)) \leq \sum_{j=1}^M u_j \log(-\mu) = \log(-\mu),$$

709 and so

$$\phi_{\ell}(\mu) \leq \log \left(\mu \sum_{j=1}^M \frac{u_j}{\mu + \ell_j} \right) = \log \left(\sum_{j=1}^M \frac{u_j}{1 + \frac{\ell_j}{\mu}} \right) \rightarrow 0 \quad \text{as } \mu \rightarrow -\infty.$$

710 This completes the first step, namely showing that there does indeed exist a unique solution $\mu^*(\tilde{\mathbf{u}})$ in
 711 the set $(-\ell_1, \infty)$ to the first equation in line (14).

712 We now turn to the second step, namely showing that this solution corresponds to the maximum total
 713 risk. Given a value of the Lagrange multiplier μ , substitution into Equation 13 gives

$$e^{t_j}(\mu) = \frac{\frac{u_j}{\mu + \ell_j}}{\sum_{k=1}^M \frac{u_k}{\mu + \ell_k}}$$

714 and therefore total risk

$$R(\mu) = \frac{\sum_{j=1}^M \frac{u_j \ell_j}{\mu + \ell_j}}{\sum_{k=1}^M \frac{u_k}{\mu + \ell_k}}.$$

715 To prove that the solution $\mu^*(\tilde{\mathbf{u}}) \in (-\infty, -\max_j \ell_j)$ is the solution to the first equation in line (14)
 716 that maximises R , it suffices to show that $R(\mu) \rightarrow \sum_{j=1}^M u_j \ell_j$ as $|\mu| \rightarrow \infty$ and $R'(\mu) \geq 0$ for all
 717 $\mu \in (-\infty, -\max_j \ell_j) \cup (-\min_j \ell_j, \infty)$, so that

$$\inf_{\mu \in (-\infty, -\max_j \ell_j)} R(\mu) \geq \sup_{\mu \in (-\min_j \ell_j, \infty)} R(\mu).$$

718 This suffices as we have already proved that $\mu^*(\tilde{\mathbf{u}})$ is the only solution in $(-\infty, -\max_j \ell_j)$ to the
 719 first equation in line (14), and that no solutions exists in the set $[-\max_j \ell_j, -\min_j \ell_j]$.

720 The limit can be easily evaluated by first rewriting $R(\mu)$ and then taking the limit as $|\mu| \rightarrow \infty$ as
 721 follows

$$R(\mu) = \frac{\sum_{j=1}^M \frac{u_j \ell_j}{1 + \frac{\ell_j}{\mu}}}{\sum_{k=1}^M \frac{u_k}{1 + \frac{\ell_k}{\mu}}} \rightarrow \frac{\sum_{j=1}^M u_j \ell_j}{\sum_{k=1}^M u_k} = \sum_{j=1}^M u_j \ell_j.$$

722 To show that $R'(\mu) \geq 0$, let $\ell_{(j)}$ denote the j^{th} smallest component of ℓ (breaking ties arbitrarily),
 723 so that $\ell_{(1)} \leq \dots \leq \ell_{(M)}$, and use the quotient rule to see that

$$\begin{aligned} R'(\mu) \geq 0 &\iff \frac{\left(\sum_{k=1}^M \frac{u_k}{\mu + \ell_k} \right) \left(\sum_{j=1}^M \frac{-u_j \ell_j}{(\mu + \ell_j)^2} \right) - \left(\sum_{j=1}^M \frac{u_j \ell_j}{\mu + \ell_j} \right) \left(\sum_{k=1}^M \frac{-u_k}{(\mu + \ell_k)^2} \right)}{\left(\sum_{p=1}^M \frac{u_p}{\mu + \ell_p} \right)^2} \geq 0 \\ &\iff \sum_{j=1}^M \sum_{k=1}^M \frac{u_j u_k \ell_j}{(\mu + \ell_j)(\mu + \ell_k)} \left(\frac{1}{\mu + \ell_k} - \frac{1}{\mu + \ell_j} \right) \geq 0 \end{aligned}$$

$$\begin{aligned} &\Leftrightarrow \sum_{\substack{j,k \in [M] \\ k < j}} \frac{u_j u_k \ell_{(j)}}{(\mu + \ell_{(j)})(\mu + \ell_{(k)})} \left(\frac{1}{\mu + \ell_{(k)}} - \frac{1}{\mu + \ell_{(j)}} \right) \\ &\quad + \sum_{\substack{j,k \in [M] \\ k > j}} \frac{u_j u_k \ell_{(j)}}{(\mu + \ell_{(j)})(\mu + \ell_{(k)})} \left(\frac{1}{\mu + \ell_{(k)}} - \frac{1}{\mu + \ell_{(j)}} \right) \geq 0, \end{aligned}$$

724 where in the final line we have dropped the summands where $k = j$ since they equal zero as the terms
725 in the bracket cancel. This final inequality holds since the first sum can be bounded below by the
726 negative of the second sum as follows

$$\begin{aligned} &\sum_{\substack{j,k \in [M] \\ k < j}} \frac{u_j u_k \ell_{(j)}}{(\mu + \ell_{(j)})(\mu + \ell_{(k)})} \left(\frac{1}{\mu + \ell_{(k)}} - \frac{1}{\mu + \ell_{(j)}} \right) \\ &\geq \sum_{\substack{j,k \in [M] \\ k < j}} \frac{u_j u_k \ell_{(k)}}{(\mu + \ell_{(j)})(\mu + \ell_{(k)})} \left(\frac{1}{\mu + \ell_{(k)}} - \frac{1}{\mu + \ell_{(j)}} \right) \quad (\text{since } \ell_{(k)} \leq \ell_{(j)} \text{ for } k < j) \\ &= \sum_{\substack{j,k \in [M] \\ k > j}} \frac{u_k u_j \ell_{(j)}}{(\mu + \ell_{(k)})(\mu + \ell_{(j)})} \left(\frac{1}{\mu + \ell_{(j)}} - \frac{1}{\mu + \ell_{(k)}} \right) \quad (\text{swapping dummy variables } j, k). \end{aligned}$$

727 We now turn to finding the partial derivatives of $F(\mathbf{t}^*(\tilde{\mathbf{u}}))$ with respect the \tilde{u}_j , which in turn will
728 allow us to find the partial derivatives of $\text{kl}_\ell^{-1}(\mathbf{u}|c)$. Let $\nabla_{\tilde{\mathbf{u}}}$ denote the gradient operator with respect
729 to $\tilde{\mathbf{u}}$. Then the quantity we are after is $\nabla_{\tilde{\mathbf{u}}} F(\mathbf{t}^*(\tilde{\mathbf{u}})) \in \mathbb{R}^{M+1}$, the j 'th component of which is

$$(\nabla_{\tilde{\mathbf{u}}} F(\mathbf{t}^*(\tilde{\mathbf{u}})))_j = \sum_{k=1}^{M+1} \frac{\partial F}{\partial t_k}(\mathbf{t}^*(\tilde{\mathbf{u}})) \frac{\partial t_k^*}{\partial \tilde{u}_j}(\tilde{\mathbf{u}}) = F_{\mathbf{t}}(\mathbf{t}^*(\tilde{\mathbf{u}})) \cdot \frac{\partial \mathbf{t}^*}{\partial \tilde{u}_j}(\tilde{\mathbf{u}}) \in \mathbb{R}.$$

730 Thus the full gradient vector is

$$\nabla_{\tilde{\mathbf{u}}} F(\mathbf{t}^*(\tilde{\mathbf{u}})) = F_{\mathbf{t}}(\mathbf{t}^*(\tilde{\mathbf{u}})) \nabla_{\tilde{\mathbf{u}}} \mathbf{t}^*(\tilde{\mathbf{u}}), \quad (15)$$

731 where $\nabla_{\tilde{\mathbf{u}}} \mathbf{t}^*(\tilde{\mathbf{u}})$ is the $M \times (M + 1)$ matrix given by

$$(\nabla_{\tilde{\mathbf{u}}} \mathbf{t}^*(\tilde{\mathbf{u}}))_{j,k} = \frac{\partial t_k^*}{\partial \tilde{u}_j}(\tilde{\mathbf{u}}).$$

732 Finding an expression for this matrix is difficult. Fortunately we can avoid needing to by using a trick
733 from mathematical economics referred to as the envelope theorem, as we now show.

734 First, note that since, for all $\tilde{\mathbf{u}}$, the constraints $g = h = 0$ are satisfied by $\mathbf{t}^*(\tilde{\mathbf{u}})$, we have the identities

$$g(\mathbf{t}^*(\tilde{\mathbf{u}}), \tilde{\mathbf{u}}) \equiv 0 \quad \text{and} \quad h(\mathbf{t}^*(\tilde{\mathbf{u}})) \equiv 0.$$

735 Differentiating these identities with respect to \tilde{u}_j then yields

$$g_{\mathbf{t}}(\mathbf{t}^*(\tilde{\mathbf{u}}), \tilde{\mathbf{u}}) \cdot \frac{\partial \mathbf{t}^*}{\partial \tilde{u}_j}(\tilde{\mathbf{u}}) + g_{\tilde{\mathbf{u}}_j}(\mathbf{t}^*(\tilde{\mathbf{u}}), \tilde{\mathbf{u}}) \equiv 0 \quad \text{and} \quad h_{\mathbf{t}}(\mathbf{t}^*(\tilde{\mathbf{u}})) \cdot \frac{\partial \mathbf{t}^*}{\partial \tilde{u}_j}(\tilde{\mathbf{u}}) \equiv 0.$$

736 As before, we can write these $M + 1$ pairs of equations as the following pair of matrix equations

$$g_{\mathbf{t}}(\mathbf{t}^*(\tilde{\mathbf{u}}), \tilde{\mathbf{u}}) \nabla_{\tilde{\mathbf{u}}} \mathbf{t}^*(\tilde{\mathbf{u}}) + g_{\tilde{\mathbf{u}}}(\mathbf{t}^*(\tilde{\mathbf{u}}), \tilde{\mathbf{u}}) \equiv \mathbf{0} \quad \text{and} \quad h_{\mathbf{t}}(\mathbf{t}^*(\tilde{\mathbf{u}})) \nabla_{\tilde{\mathbf{u}}} \mathbf{t}^*(\tilde{\mathbf{u}}) \equiv \mathbf{0}.$$

737 Multiplying these identities by $\lambda^*(\tilde{\mathbf{u}})$ and $\mu^*(\tilde{\mathbf{u}})$ respectively, and combining with equation (15),
738 yields

$$\begin{aligned} \nabla_{\tilde{\mathbf{u}}} F(\mathbf{t}^*(\tilde{\mathbf{u}})) &= \left(F_{\mathbf{t}}(\mathbf{t}^*(\tilde{\mathbf{u}})) + \lambda^*(\tilde{\mathbf{u}}) g_{\mathbf{t}}(\mathbf{t}^*(\tilde{\mathbf{u}}), \tilde{\mathbf{u}}) + \mu^*(\tilde{\mathbf{u}}) h_{\mathbf{t}}(\mathbf{t}^*(\tilde{\mathbf{u}})) \right) \nabla_{\tilde{\mathbf{u}}} \mathbf{t}^*(\tilde{\mathbf{u}}) \\ &\quad + \lambda^*(\tilde{\mathbf{u}}) g_{\tilde{\mathbf{u}}}(\mathbf{t}^*(\tilde{\mathbf{u}}), \tilde{\mathbf{u}}) \\ &= \lambda^*(\tilde{\mathbf{u}}) g_{\tilde{\mathbf{u}}}(\mathbf{t}^*(\tilde{\mathbf{u}}), \tilde{\mathbf{u}}), \end{aligned}$$

739 where the final equality comes from noting that the terms in the large bracket vanish due to equation
 740 (12). Recalling the expression for $g_{\tilde{\mathbf{u}}}(\mathbf{t}; \tilde{\mathbf{u}})$ given by Equation 11 and that $\mathbf{v}^*(\tilde{\mathbf{u}}) = \exp(\mathbf{t}^*(\tilde{\mathbf{u}}))$ we
 741 obtain

$$\begin{aligned} \nabla_{\tilde{\mathbf{u}}} F(\mathbf{t}^*(\tilde{\mathbf{u}})) &= \lambda^*(\tilde{\mathbf{u}}) \left(1 - \mathbf{t}^*(\tilde{\mathbf{u}})_1 + \log u_1, \dots, 1 - \mathbf{t}^*(\tilde{\mathbf{u}})_M + \log u_M, -1 \right) \\ &= \lambda^*(\tilde{\mathbf{u}}) \left(1 + \log \frac{u_1}{\mathbf{v}^*(\tilde{\mathbf{u}})_1}, \dots, 1 + \log \frac{u_M}{\mathbf{v}^*(\tilde{\mathbf{u}})_M}, -1 \right) \end{aligned}$$

742 Finally, recalling Equivalence (10), namely $\nabla_{\tilde{\mathbf{u}}} f_{\ell}^*(\tilde{\mathbf{u}}) \equiv \nabla_{\tilde{\mathbf{u}}} F(\mathbf{t}^*(\tilde{\mathbf{u}}))$, we see that the above
 743 expression gives the derivatives $\frac{\partial f_{\ell}^*}{\partial u_j}(\tilde{\mathbf{u}})$ and $\frac{\partial f_{\ell}^*}{\partial c}(\tilde{\mathbf{u}})$ stated in the proposition, thus completing the
 744 proof. □