

A Minimization of L_2 WD between Univariate Gaussian and Non-Gaussian Distributions

In this section, we derive the formulas of the optimal μ^* and σ^* for the L_2 WD, *i.e.*, Eqn. (5). Recall the optimization problem: we use a univariate Gaussian distribution $\mathcal{N}(f|\mu, \sigma^2)$ to approximate a univariate non-Gaussian distribution $q(f)$ by minimizing the L_2 WD between them:

$$\min_{\mu, \sigma} \mathbb{W}_2^2(q, \mathcal{N}) = \min_{\mu, \sigma} \int_0^1 \left| F_q^{-1}(y) - \mu - \sqrt{2}\sigma \operatorname{erf}^{-1}(2y - 1) \right|^2 dy,$$

where F_q^{-1} is the quantile function of the non-Gaussian distribution q , namely the pseudoinverse function of the corresponding cumulative distribution function F_q defined in Proposition 1.

To solve this problem, we first calculate derivatives about μ and σ :

$$\begin{aligned} \frac{\partial \mathbb{W}_2^2}{\partial \mu} &= -2 \int_0^1 F_q^{-1}(y) - \mu - \sqrt{2}\sigma \operatorname{erf}^{-1}(2y - 1) dy, \\ \frac{\partial \mathbb{W}_2^2}{\partial \sigma} &= -2 \int_0^1 (F_q^{-1}(y) - \mu - \sqrt{2}\sigma \operatorname{erf}^{-1}(2y - 1)) \sqrt{2} \operatorname{erf}^{-1}(2y - 1) dy. \end{aligned}$$

Then, by zeroing derivatives, we obtain the optimal parameters:

$$\begin{aligned} \mu^* &= \int_0^1 F_q^{-1}(y) - \sqrt{2}\sigma \operatorname{erf}^{-1}(2y - 1) dy \\ &= \int_{-\infty}^{\infty} xq(x) dx - \frac{\sqrt{2}}{2}\sigma \int_{-1}^1 \operatorname{erf}^{-1}(y) dy \\ &= \mu_q - \sqrt{2}\sigma \int_{-\infty}^{\infty} x\mathcal{N}(x|0, 1/2) dx \\ &= \mu_q, \\ \sigma^* &= \sqrt{2} \int_0^1 (F_q^{-1}(y) - \mu) \operatorname{erf}^{-1}(2y - 1) dy / \int_0^1 2(\operatorname{erf}^{-1})^2(2y - 1) dy \\ &= \sqrt{2} \int_0^1 F_q^{-1}(y) \operatorname{erf}^{-1}(2y - 1) dy / \underbrace{\int_{-\infty}^{\infty} 2x^2 \mathcal{N}(x|0, 1/2) dx}_{=1} \\ &= \sqrt{2} \int_0^1 F_q^{-1}(y) \operatorname{erf}^{-1}(2y - 1) dy \\ &= \sqrt{2} \int_{-\infty}^{\infty} f \operatorname{erf}^{-1}(2F_q(f) - 1) dF_q(f) \\ &= -\sqrt{\frac{1}{2\pi}} \int_{-\infty}^{\infty} f d e^{-[\operatorname{erf}^{-1}(2F_q(f)-1)]^2} \\ &= 0 + \sqrt{\frac{1}{2\pi}} \int_{-\infty}^{\infty} e^{-[\operatorname{erf}^{-1}(2F_q(f)-1)]^2} df. \end{aligned} \tag{8}$$

B Minimization of L_p WD between Univariate Gaussian and Non-Gaussian Distributions

In this section, we describe a gradient descent approach to minimizing an L_p WD, for $p \neq 2$, in order to handle cases with no analytical expressions for the optimal parameters. Our goal is to use a univariate Gaussian distribution $\mathcal{N}(f|\mu, \sigma^2)$ to approximate a univariate non-Gaussian distribution $q(f)$. Specifically, we seek the minimiser in μ and σ of $\mathbb{W}_p^p(q, \mathcal{N})$; the derivatives of the objective

function about μ and σ are:

$$\begin{aligned}\partial_\mu \mathbf{W}_p^p &= -p \int_0^1 |\varepsilon(y)|^{p-1} \text{sgn}(\varepsilon(y)) \, dy = -p \int_{-\infty}^{\infty} |\eta(x)|^{p-1} \text{sgn}(\eta(x)) q(x) \, dx, \\ \partial_\sigma \mathbf{W}_p^p &= -p \int_0^1 |\varepsilon(y)|^{p-1} \text{sgn}(\varepsilon(y)) \text{erf}^{-1}(2y-1) \, dy = -p \int_{-\infty}^{\infty} |\eta(x)|^{p-1} \text{sgn}(\eta(x)) \text{erf}^{-1}(2F_q(x)-1) q(x) \, dx.\end{aligned}$$

where for simplification, we define $\varepsilon(y) = F_q^{-1}(y) - \mu - \sqrt{2}\sigma \text{erf}^{-1}(2y-1)$ and $\eta(x) = x - \mu - \sqrt{2}\sigma \text{erf}^{-1}(2F_q(x)-1)$, with F_q and F_q^{-1} being the CDF and the quantile function of q . Note the derivatives have no analytical expressions. However, if the CDF F_q is available, we can use the standard numerical integration routines; otherwise, we resort to Monte Carlo sampling. In the framework of EP or QP, $q(x) \propto q^{\setminus i}(x)p(y_i|x)$ and $q^{\setminus i}$ is Gaussian, so we may draw samples from a Gaussian proposal distribution to obtain a simple Monte Carlo method.

C Computations for Different Likelihoods

Given the likelihood $p(y|f)$ and the cavity distribution $q^{\setminus i}(f) = \mathcal{N}(f|\mu, \sigma^2)$, a stable way to compute the mean and the variance of the tilted distribution $\tilde{q}(f) = p(y|f)q^{\setminus i}(f)/Z$ where the normalizer $Z = \int_{-\infty}^{\infty} p(y|f)q^{\setminus i}(f) \, df$, can be found in the software manual of Rasmussen and Williams [47]. We present the key formulae below, for use in subsequent derivations:

$$\begin{aligned}\partial_\mu Z &= \int_{-\infty}^{\infty} \frac{f-\mu}{\sigma^2} p(y|f) \mathcal{N}(f|\mu, \sigma^2) \, df \\ \frac{\partial_\mu Z}{Z} &= \frac{1}{\sigma^2} \int_{-\infty}^{\infty} f \frac{p(y|f) \mathcal{N}(f|\mu, \sigma^2)}{Z} \, df - \frac{\mu}{\sigma^2} \int_{-\infty}^{\infty} \frac{p(y|f) \mathcal{N}(f|\mu, \sigma^2)}{Z} \, dy \\ \frac{\partial_\mu Z}{Z} &= \frac{1}{\sigma^2} \mu_{\tilde{q}} - \frac{\mu}{\sigma^2} \\ \implies \mu_{\tilde{q}} &= \frac{\sigma^2 \partial_\mu Z}{Z} + \mu = \sigma^2 \partial_\mu \log Z + \mu, \\ \partial_\mu^2 Z &= \int_{-\infty}^{\infty} -\frac{1}{\sigma^2} p(y|f) \mathcal{N}(f|\mu, \sigma^2) + \left(\frac{f-\mu}{\sigma^2} \right)^2 p(y|f) \mathcal{N}(f|\mu, \sigma^2) \, df \\ \frac{\partial_\mu^2 Z}{Z} &= \int_{-\infty}^{\infty} \left(-\frac{1}{\sigma^2} + \frac{\mu^2}{\sigma^4} + \frac{f^2}{\sigma^4} - \frac{2\mu f}{\sigma^4} \right) \frac{p(y|f) \mathcal{N}(f|\mu, \sigma^2)}{Z} \, df \\ \frac{\partial_\mu^2 Z}{Z} &= -\frac{1}{\sigma^2} + \frac{\mu^2}{\sigma^4} + \frac{1}{\sigma^4} (\sigma_{\tilde{q}}^2 + \mu_{\tilde{q}}^2) - \frac{2\mu}{\sigma^4} \mu_{\tilde{q}} \\ \frac{\partial_\mu^2 Z}{Z} &= -\frac{1}{\sigma^2} + \frac{\sigma_{\tilde{q}}^2}{\sigma^4} + \frac{(\mu - \mu_{\tilde{q}})^2}{\sigma^4} = -\frac{1}{\sigma^2} + \frac{\sigma_{\tilde{q}}^2}{\sigma^4} + \left(\frac{\partial_\mu Z}{Z} \right)^2 \\ \implies \sigma_{\tilde{q}}^2 &= \sigma^4 \left[\frac{\partial_\mu^2 Z}{Z} - \left(\frac{\partial_\mu Z}{Z} \right)^2 \right] + \sigma^2 = \sigma^4 \partial_\mu^2 \log Z + \sigma^2.\end{aligned}$$

C.1 Probit Likelihood for Binary Classification

For the binary classification with labels $y \in \{-1, 1\}$, the PDF of the tilted distribution $\tilde{q}(f)$ with the probit likelihood is provided by Rasmussen and Williams [47]:

$$\tilde{q}(f) = Z^{-1} \Phi(fy) \mathcal{N}(f|\mu, \sigma^2), \quad Z = \Phi(z), \quad z = \frac{\mu}{y\sqrt{1+\sigma^2}},$$

and the mean estimate also has a closed form expression:

$$\mu^* = \mu_{\tilde{q}} = \mu + \frac{\sigma^2 \mathcal{N}(z)}{\Phi(z) y \sqrt{1+\sigma^2}}.$$

As per Equation (5), the computation of the optimal σ^* requires the CDF of \tilde{q} , denoted as $F_{\tilde{q}}$. For positive $y > 0$, the CDF is derived as

$$\begin{aligned}
F_{\tilde{q}, y > 0}(x) &= Z^{-1} \int_{-\infty}^x \Phi(fy) \mathcal{N}(f|\mu, \sigma^2) \, df \\
&= \frac{Z^{-1}}{2\pi\sigma y} \int_{-\infty}^{\mu} \int_{-\infty}^{x-\mu} \exp\left(-\frac{1}{2} \begin{bmatrix} w \\ f \end{bmatrix}^T \begin{bmatrix} v^{-2} + \sigma^{-2} & v^{-2} \\ v^{-2} & v^{-2} \end{bmatrix} \begin{bmatrix} w \\ f \end{bmatrix}\right) \, dw \, df \\
&= Z^{-1} \int_{-\infty}^k \int_{-\infty}^h \mathcal{N}\left(\begin{bmatrix} w \\ f \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} 1 & -\rho \\ -\rho & 1 \end{bmatrix}\right) \, dw \, df \\
&\stackrel{(a)}{=} Z^{-1} \left[\frac{1}{2} \Phi(h) - T\left(h, \frac{k + \rho h}{h\sqrt{1-\rho^2}}\right) + \frac{1}{2} \Phi(k) - T\left(k, \frac{h + \rho k}{k\sqrt{1-\rho^2}}\right) + \eta \right] \\
&k = \frac{\mu}{\sqrt{\sigma^2 + 1}}, \quad h = \frac{x - \mu}{\sigma}, \quad \rho = \frac{1}{\sqrt{1 + 1/\sigma^2}}, \quad x \neq \mu, \quad \mu \neq 0,
\end{aligned}$$

where the step (a) is obtained by exploiting the work of Owen [45] and $T(\cdot, \cdot)$ is the Owen's T function:

$$T(h, a) = \frac{1}{2\pi} \int_0^a \frac{\exp[-(1+x^2)h^2/2]}{1+x^2} \, dx,$$

and η is defined as

$$\eta = \begin{cases} 0 & hk > 0 \text{ or } (hk = 0 \text{ and } h + k \geq 0), \\ -0.5 & \text{otherwise.} \end{cases}$$

Similarly, for $y < 0$, the CDF is

$$F_{\tilde{q}, y < 0}(x) = Z^{-1} \left[\frac{1}{2} \Phi(h) + T\left(h, \frac{k + \rho h}{h\sqrt{1-\rho^2}}\right) - \frac{1}{2} \Phi(k) + T\left(k, \frac{h + \rho k}{k\sqrt{1-\rho^2}}\right) - \eta \right].$$

Summarizing the two cases, we get the closed form expression of $F_{\tilde{q}}$:

$$\begin{aligned}
F_{\tilde{q}}(x) &= Z^{-1} \left[\frac{1}{2} \Phi(h) - yT\left(h, \frac{k + \rho h}{h\sqrt{1-\rho^2}}\right) + \frac{y}{2} \Phi(k) - yT\left(k, \frac{h + \rho k}{k\sqrt{1-\rho^2}}\right) + y\eta \right] \\
&= Z^{-1} \left[\frac{1}{2} \Phi(h) - yT\left(h, \frac{k}{h\sqrt{1-\rho^2}} + \sigma\right) + \frac{y}{2} \Phi(k) - yT\left(k, \frac{h}{k\sqrt{1-\rho^2}} + \sigma\right) + y\eta \right].
\end{aligned}$$

Provided the above, the optimal σ^* can be computed by numerical integration of Eqn (8). For special cases, we provide additional formulas:

$$\begin{aligned}
(1) \quad x = \mu, \quad \mu \neq 0 : F_{\tilde{q}}(x) &= Z^{-1} \left[\frac{1}{4} - \frac{y \operatorname{sign}(k)}{4} + \frac{y}{2} \Phi(k) - yT(k, \sigma) + y\eta \right]; \\
(2) \quad x \neq \mu, \quad \mu = 0 : F_{\tilde{q}}(x) &= 2 \left[\frac{1}{2} \Phi(h) - yT(h, \sigma) + \frac{y}{4} - \frac{y \operatorname{sign}(h)}{4} + y\eta \right]; \\
(3) \quad x = \mu, \quad \mu = 0 : F_{\tilde{q}}(x) &= \frac{1}{2} - \frac{y}{\pi} \arctan(\sigma).
\end{aligned}$$

C.2 Square Link Function for Poisson Regression

Consider Poisson regression, which uses the Poisson likelihood $p(y|g) = g^y \exp(-g)/y!$ to model count data $y \in \mathbb{N}$, with the square link function $g(f) = f^2$ [56, 15]. We use the square link function because it is more mathematically convenient than the exponential function. Given the cavity distribution $q^{\setminus i}(f) = \mathcal{N}(f|\mu, \sigma^2)$, we want the tilted distribution $\tilde{q}(f) = q^{\setminus i}(f)p(y|g(f))/Z$ where the normalizer Z is derived as:

$$Z = \int_{-\infty}^{\infty} q^{\setminus i}(f)p(y|g) \, df$$

$$\begin{aligned}
&= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(f-\mu)^2}{2\sigma^2}\right) f^{2y} \exp(-f^2)/y! \, df \\
&\stackrel{(a)}{=} \frac{1}{\sqrt{2\pi\sigma^2}y! \exp(\mu^2/(1+2\sigma^2))} \int_{-\infty}^{\infty} f^{2y} \exp\left(-\frac{(f-\mu/(1+2\sigma^2))^2}{2\sigma^2/(1+2\sigma^2)}\right) \, df \\
&\stackrel{(b)}{=} \frac{\left(\frac{2\sigma^2}{1+2\sigma^2}\right)^{y+\frac{1}{2}}}{\sqrt{2\pi\sigma^2}y! \exp(\mu^2/(1+2\sigma^2))} \Gamma\left(y+\frac{1}{2}\right) {}_1F_1\left(-y; \frac{1}{2}; -\frac{\mu^2}{2\sigma^2(1+2\sigma^2)}\right) \\
&= \frac{\alpha^{y+\frac{1}{2}}}{\sqrt{2\pi\sigma^2}y! \exp(h)} \Gamma\left(y+\frac{1}{2}\right) {}_1F_1\left(-y; \frac{1}{2}; -\frac{h}{2\sigma^2}\right), \\
\alpha &= \frac{2\sigma^2}{1+2\sigma^2}, \quad h = \frac{\mu^2}{1+2\sigma^2} \tag{9}
\end{aligned}$$

where the step (a) rewrites the product of two exponential functions into the form of the Gaussian distribution, (b) is achieved through Mathematica [59], $\Gamma(\cdot)$ is the Gamma function and ${}_1F_1\left(-y; \frac{1}{2}; -\frac{h}{2\sigma^2}\right)$ is the confluent hypergeometric function of the first kind. Furthermore, we compute the first derivative of $\log Z$ w.r.t. μ and then the mean of the tilted distribution:

$$\begin{aligned}
\partial_{\mu} \log Z &= \left(\frac{y {}_1F_1\left(-y+1; \frac{3}{2}; -\frac{h}{2\sigma^2}\right)}{\sigma^2 {}_1F_1\left(-y; \frac{1}{2}; -\frac{h}{2\sigma^2}\right)} - 1 \right) \frac{2\mu}{1+2\sigma^2} \\
&\Rightarrow \mu_{\tilde{q}} = \sigma^2 \partial_{\mu} \log Z + \mu. \\
\partial_{\mu}^2 \log Z &= \left(\frac{y {}_1F_1\left(-y+1; \frac{3}{2}; -\frac{h}{2\sigma^2}\right)}{\sigma^2 {}_1F_1\left(-y; \frac{1}{2}; -\frac{h}{2\sigma^2}\right)} - 1 \right) \frac{2}{1+2\sigma^2} - \\
&\quad \left(\frac{2(1-y) {}_1F_1\left(-y+2; \frac{5}{2}; -\frac{h}{2\sigma^2}\right)}{3 {}_1F_1\left(-y; \frac{1}{2}; -\frac{h}{2\sigma^2}\right)} + \frac{2y {}_1F_1\left(-y+1; \frac{3}{2}; -\frac{h}{2\sigma^2}\right)^2}{{}_1F_1\left(-y; \frac{1}{2}; -\frac{h}{2\sigma^2}\right)^2} \right) \frac{2\mu^2 y}{\sigma^4(1+2\sigma^2)^2} \\
&\Rightarrow \sigma_{\tilde{q}}^2 = \sigma^4 \partial_{\mu}^2 \log Z + \sigma^2
\end{aligned}$$

Finally, we derive the CDF of the tilted distribution \tilde{q} by using the binomial theorem:

$$\begin{aligned}
F_{\tilde{q}}(x) &= Z^{-1} \int_{-\infty}^x p(y|g) \mathcal{N}(f|\mu, \sigma^2) \, df \\
&\stackrel{(a)}{=} A \int_{-\infty}^x f^{2y} \exp\left(-\frac{(f-\mu/(1+2\sigma^2))^2}{2\sigma^2/(1+2\sigma^2)}\right) \, df \\
&= A \int_{-\infty}^{x-\frac{\mu}{1+2\sigma^2}} \left(f + \frac{\mu}{1+2\sigma^2}\right)^{2y} \exp\left(-\frac{f^2}{2\sigma^2/(1+2\sigma^2)}\right) \, df \\
&\stackrel{(b)}{=} A \int_{-\infty}^{x-\beta} \left[\sum_{k=0}^{2y} \binom{2y}{k} f^k \beta^{2y-k} \right] \exp\left(-\frac{f^2}{\alpha}\right) \, df \\
&= A \sum_{k=0}^{2y} \binom{2y}{k} \beta^{2y-k} \left[\int_{-\infty}^0 f^k \exp\left(-\frac{f^2}{\alpha}\right) \, df + \int_0^{x-\beta} f^k \exp\left(-\frac{f^2}{\alpha}\right) \, df \right] \\
&\stackrel{(c)}{=} \frac{A}{2} \sum_{k=0}^{2y} \binom{2y}{k} \beta^{2y-k} \alpha^{\frac{k+1}{2}} \left[(-1)^k \Gamma\left(\frac{k+1}{2}\right) + \operatorname{sgn}(x-\beta)^{k+1} \left(\Gamma\left(\frac{k+1}{2}\right) - \Gamma\left(\frac{k+1}{2}, \frac{(x-\beta)^2}{\alpha}\right) \right) \right] \\
A &= \frac{Z^{-1}}{\sqrt{2\pi\sigma^2}y! \exp(\mu^2/(1+2\sigma^2))} = \left[\alpha^{y+\frac{1}{2}} \Gamma\left(y+\frac{1}{2}\right) {}_1F_1\left(-y; \frac{1}{2}; -\frac{h}{2\sigma^2}\right) \right]^{-1}, \quad \beta = \frac{\mu}{1+2\sigma^2},
\end{aligned}$$

where the step (a) has been derived in (a) of Eqn. (9), (b) applies the binomial theorem and (c) is obtained through Mathematica [59]. And, the function $\Gamma(a, z) = \int_z^{\infty} t^{a-1} e^{-t} \, dt$ is the upper

incomplete gamma function and $\text{sgn}(x)$ is the sign function, equaling 1 when $x > 0$, 0 when $x = 0$ and -1 when $x < 0$.

D Proof of Convexity

Theorem Given two probability measures in $\mathcal{M}_+^1(\mathbb{R})$: a Gaussian $\mathcal{N}(\mu, \sigma^2)$ with mean μ and standard deviation $\sigma > 0$, and an arbitrary measure \tilde{q} , the L_p WD $\mathbb{W}_p^p(\tilde{q}, \mathcal{N})$ is strictly convex about μ and σ .

Proof. Let $F_{\tilde{q}}^{-1}(y)$ and $F_{\mathcal{N}}^{-1}(y) = \mu + \sqrt{2}\sigma\text{erf}^{-1}(2y - 1)$, $y \in [0, 1]$, be the quantile functions of \tilde{q} and the Gaussian \mathcal{N} , where erf is the error function. Then, we consider two distinct Gaussian measures $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$ and a convex combination w.r.t. their parameters $\mathcal{N}(a_1\mu_1 + a_2\mu_2, (a_1\sigma_1 + a_2\sigma_2)^2)$ with $a_1, a_2 \in \mathbb{R}_+$ and $a_1 + a_2 = 1$. Given the above, we further define $\varepsilon_k(y) = F_{\tilde{q}}^{-1}(y) - \mu_k - \sigma_k\sqrt{2}\text{erf}^{-1}(2y - 1)$, $k = 1, 2$, for notational simplification, and derive the convexity as:

$$\begin{aligned} \mathbb{W}_p^p(\tilde{q}, \mathcal{N}(a_1\mu_1 + a_2\mu_2, (a_1\sigma_1 + a_2\sigma_2)^2)) &\stackrel{(a)}{=} \int_0^1 |a_1\varepsilon_1(y) + a_2\varepsilon_2(y)|^p \, dy \stackrel{(b)}{\leq} \int_0^1 (a_1|\varepsilon_1(y)| + \\ &a_2|\varepsilon_2(y)|)^p \, dy \stackrel{(c)}{\leq} a_1\mathbb{W}_p^p(\tilde{q}, \mathcal{N}(\mu_1, \sigma_1^2)) + a_2\mathbb{W}_p^p(\tilde{q}, \mathcal{N}(\mu_2, \sigma_2^2)), \end{aligned}$$

where steps (a), (b) and (c) are obtained by applying Proposition 1, non-negativity of the absolute value, and the convexity of $f(x) = x^p$, $p \geq 1$, over \mathbb{R}_+ respectively. The equality at (b) holds iff $\varepsilon_k(y) \geq 0$, $k = 1, 2$, $\forall y \in [0, 1]$, and (c)'s equality holds iff $|\varepsilon_1(y)| = |\varepsilon_2(y)|$, $\forall y \in [0, 1]$. These two conditions for equality can't be attained simultaneously as otherwise it would contradict that $\mathcal{N}(\mu_1, \sigma_1^2)$ is different from $\mathcal{N}(\mu_2, \sigma_2^2)$. So, $\mathbb{W}_p^p(\tilde{q}, \mathcal{N})$, $p \geq 1$, is strictly convex about μ and σ . \square

E Proof of Variance Difference

Theorem The variance of the Gaussian approximation to a univariate tilted distribution $\tilde{q}(f)$ as estimated by QP and EP satisfy $\sigma_{\text{QP}}^2 \leq \sigma_{\text{EP}}^2$.

Proof. Let $\mathcal{N}(\mu_{\text{QP}}, \sigma_{\text{QP}}^2)$ be the optimal Gaussian in QP. As per Proposition 1, we reformulate the L_2 WD based projection $\mathbb{W}_2^2(\tilde{q}, \mathcal{N}(\mu_{\text{QP}}, \sigma_{\text{QP}}^2))$ w.r.t. quantile functions:

$$\begin{aligned} \mathbb{W}_2^2(\tilde{q}, \mathcal{N}(\mu_{\text{QP}}, \sigma_{\text{QP}}^2)) &= \int_0^1 |F_{\tilde{q}}^{-1}(y) - \mu_{\text{QP}} - \sqrt{2}\sigma_{\text{QP}}\text{erf}^{-1}(2y - 1)|^2 \, dy = \int_0^1 \underbrace{(F_{\tilde{q}}^{-1}(y) - \mu_{\text{QP}})^2}_{\sigma_{\text{EP}}^2} \\ &+ \underbrace{(\sqrt{2}\sigma_{\text{QP}}\text{erf}^{-1}(2y - 1))^2}_{\sigma_{\text{QP}}^2} - 2 \underbrace{(F_{\tilde{q}}^{-1}(y) - \mu_{\text{QP}})\sqrt{2}\sigma_{\text{QP}}\text{erf}^{-1}(2y - 1)}_{(A)} \, dy = \sigma_{\text{EP}}^2 - \sigma_{\text{QP}}^2, \end{aligned}$$

where for (A), we used $\int \mu_{\text{QP}}\sigma_{\text{QP}}\text{erf}^{-1}(2y - 1) \, dy = 0$ and the remaining factor can be easily shown to be equal to $2\sigma_{\text{QP}}^2$. Furthermore, due to the non-negativity of the WD, we have $\sigma_{\text{EP}}^2 \geq \sigma_{\text{QP}}^2$, and the equality holds iff \tilde{q} is Gaussian. \square

F Proof of Locality Property

Theorem Minimization of $\mathbb{W}_2^2(\tilde{q}(f), \mathcal{N}(f))$ w.r.t. $\mathcal{N}(f)$ results in $q^{\setminus i}(f_{\setminus i}|f_i) = \mathcal{N}(f_{\setminus i}|f_i)$.

Proof. We first apply the decomposition of the L_2 norm to rewriting the $\mathbb{W}_2^2(\tilde{q}(f), \mathcal{N}(f))$ as below (see detailed derivations in Appendix F.2):

$$\mathbb{W}_2^2(\tilde{q}, \mathcal{N}) = \inf_{\pi_i} \mathbb{E}_{\pi_i} \left[\|f_i - f_i'\|_2^2 + \mathbb{W}_2^2(q^{\setminus i}_{\setminus i}, \mathcal{N}_{\setminus i}) \right], \quad (10)$$

where the prime indicates that the variable is from the Gaussian \mathcal{N} , and for simplification, we use the notation π_i for the joint distribution $\pi(f_i, f_i')$ which belongs to a set of measures $U(\tilde{q}_i, \mathcal{N}_i)$. Since

$q^{\setminus i}(\mathbf{f})$ is known to be Gaussian, we define it in a partitioned form:

$$q^{\setminus i}(\mathbf{f}) \equiv \mathcal{N} \left(\begin{bmatrix} \mathbf{f}_{\setminus i} \\ f_i \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m}_{\setminus i} \\ m_i \end{bmatrix}, \begin{bmatrix} \mathbf{S}_{\setminus ii} & \mathbf{S}_{\setminus ii} \\ \mathbf{S}_{\setminus ii}^\top & S_i \end{bmatrix} \right), \quad (11)$$

and the conditional $q^{\setminus i}(\mathbf{f}_{\setminus i}|f_i)$ is expressed as:

$$q^{\setminus i}(\mathbf{f}_{\setminus i}|f_i) = \mathcal{N}(\mathbf{f}_{\setminus i}|\mathbf{m}_{\setminus i|i}, \mathbf{S}_{\setminus i|i}), \quad \mathbf{m}_{\setminus i|i} = \mathbf{m}_{\setminus i} + \mathbf{S}_{\setminus ii} \mathbf{S}_i^{-1} (f_i - m_i) \equiv \mathbf{a} f_i + \mathbf{b}, \quad (12)$$

$$\mathbf{S}_{\setminus i|i} = \mathbf{S}_{\setminus i} - \mathbf{S}_{\setminus ii} \mathbf{S}_i^{-1} \mathbf{S}_{\setminus ii}^\top.$$

We define a similar partitioned expression for the Gaussian $\mathcal{N}(\mathbf{f}')$ by adding primes to variables and parameters on the r.h.s. of Equation (11), and as a result, the conditional $\mathcal{N}(\mathbf{f}'_{\setminus i}|f'_i)$ is written as:

$$\mathcal{N}(\mathbf{f}'_{\setminus i}|f'_i) = \mathcal{N}(\mathbf{m}'_{\setminus i|i}, \mathbf{S}'_{\setminus i|i}), \quad \mathbf{m}'_{\setminus i|i} = \mathbf{m}'_{\setminus i} + \mathbf{S}'_{\setminus ii} \mathbf{S}'_i{}^{-1} (f'_i - m'_i) \equiv \mathbf{a}' f'_i + \mathbf{b}', \quad (13)$$

$$\mathbf{S}'_{\setminus i|i} = \mathbf{S}'_{\setminus i} - \mathbf{S}'_{\setminus ii} \mathbf{S}'_i{}^{-1} \mathbf{S}'_{\setminus ii}^\top. \quad (14)$$

Given the above definitions, we exploit Proposition 2 to take the means out of the L_2 WD on the r.h.s. of Equation (10):

$$\mathbf{W}_2^2(\tilde{q}, \mathcal{N}) = \inf_{\pi_i} \mathbb{E}_{\pi_i} \left[\|f_i - f'_i\|_2^2 + \|\mathbf{m}_{\setminus i|i} - \mathbf{m}'_{\setminus i|i}\|_2^2 \right] + \underbrace{\mathbf{W}_2^2 \left(\mathcal{N}(\mathbf{0}, \mathbf{S}_{\setminus i|i}), \mathcal{N}(\mathbf{0}, \mathbf{S}'_{\setminus i|i}) \right)}_{(A)}. \quad (15)$$

Minimizing this function requires optimizing m'_i , $\mathbf{m}'_{\setminus i}$, S'_i , $\mathbf{S}'_{\setminus i}$ and $\mathbf{S}'_{\setminus ii}$. As $\mathbf{S}'_{\setminus i}$ is only contained in $\mathbf{S}_{\setminus i|i}$ and isolated into the term (A), it can be optimized by simply setting

$$\mathbf{S}'_{\setminus i|i} = \mathbf{S}_{\setminus i|i} \xrightarrow{\text{Eqn. (14)}} \mathbf{S}_{\setminus i|i}^{(n)*} = \mathbf{S}_{\setminus i|i} + \mathbf{S}'_{\setminus ii} \mathbf{S}'_i{}^{-1} \mathbf{S}'_{\setminus ii}^\top. \quad (16)$$

As a result, (A) is minimized to zero. Next, we plug in expressions of $\mathbf{m}_{\setminus i|i}$ and $\mathbf{m}'_{\setminus i|i}$ (Equation (12) and Equation (13)) into optimized Equation (15):

$$\min_{\mathbf{S}'_{\setminus i}} (15) = \inf_{\pi_i} \mathbb{E}_{\pi_i} \left[\|f_i - f'_i\|_2^2 + \|\mathbf{a} f_i - \mathbf{a}' f'_i + \mathbf{b} - \mathbf{b}'\|_2^2 \right], \quad (17)$$

where $\mathbf{m}'_{\setminus i}$ is only contained by \mathbf{b}' . Thus, we can optimize it by zeroing the derivative of the above function about $\mathbf{m}'_{\setminus i}$, which results in:

$$\mathbf{b}' = \mathbf{b} + \mathbf{a} \mu_{\tilde{q}_i} - \mathbf{a}' m'_i \xrightarrow{\text{Eqn. (13)}} \mathbf{m}_{\setminus i}^{(n)*} = \mathbf{S}'_{\setminus ii} \mathbf{S}'_i{}^{-1} m'_i + \mathbf{b} + \mathbf{a} \mu_{\tilde{q}_i} - \mathbf{a}' m'_i, \quad (18)$$

where $\mu_{\tilde{q}_i}$ is the mean of $\tilde{q}(f_i)$. The minimum value of Equation (17) thereby is (see details in subsection F.3):

$$\min_{\mathbf{m}'_{\setminus i}} (17) = (1 + \mathbf{a}^\top \mathbf{a}') \mathbf{W}_2^2(\tilde{q}_i, \mathcal{N}_i) + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 + \|\mathbf{a}'\|_2^2 S'_i - \mathbf{a}^\top \mathbf{a}' \left[\sigma_{\tilde{q}_i}^2 + S'_i + (\mu_{\tilde{q}_i} - m'_i)^2 \right] \quad (19)$$

where $\sigma_{\tilde{q}_i}^2$ is the variance of $\tilde{q}(f_i)$. This function can be further simplified using the quantile based reformulation of $\mathbf{W}_2^2(\tilde{q}_i, \mathcal{N}_i)$ (see details in Appendix F.4) which results in:

$$(19) = \mathbf{W}_2^2(\tilde{q}_i, \mathcal{N}_i) + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 - \underbrace{2^{\frac{3}{2}} \mathbf{a}^\top \mathbf{a}' c_{\tilde{q}_i} S'_i{}^{\frac{1}{2}}}_{(B)} + \|\mathbf{a}'\|_2^2 S'_i. \quad (20)$$

Now, we are left with optimizing m'_i , S'_i and $\mathbf{S}'_{\setminus ii}$. To optimize $\mathbf{S}'_{\setminus ii}$, which only exists in the above term (B), we zero the derivative of (B) w.r.t. $\mathbf{S}'_{\setminus ii}$ and this yields:

$$\mathbf{a}'^* = 2^{\frac{1}{2}} (S'_i)^{-\frac{1}{2}} c_{\tilde{q}_i} \mathbf{a} \xrightarrow{\text{Eqn. (13)}} \mathbf{S}'_{\setminus ii}{}^* = (2S'_i)^{\frac{1}{2}} c_{\tilde{q}_i} \mathbf{a}, \quad (21)$$

and the minimum value of Equation (20) is

$$\min_{\mathbf{S}'_{\setminus ii}} (20) = \mathbf{W}_2^2(\tilde{q}_i, \mathcal{N}_i) + \|\mathbf{a}\|_2^2 (\sigma_{\tilde{q}_i}^2 - 2c_{\tilde{q}_i}^2). \quad (22)$$

The results of optimizing m'_i and S'_i in the above equation have already been provided in Equation (5): $m_i'^* = \mu_{\tilde{q}_i}$ and $S_i'^* = 2c_{\tilde{q}_i}^2$. By plugging them into Equation (21) and Equation (18), we have

$\mathbf{a}^* = \mathbf{a}$ and $\mathbf{b}^* = \mathbf{b}$. Finally, using Equation (16), we obtain $q^{\setminus i}(\mathbf{f}_{\setminus i}|f_i) = \mathcal{N}(\mathbf{f}_{\setminus i}|\mathbf{a}f_i + \mathbf{b}, \mathbf{S}_{\setminus i|i}) = \mathcal{N}(\mathbf{f}_{\setminus i}|\mathbf{a}'f_i + \mathbf{b}', \mathbf{S}'_{\setminus i|i}) = \mathcal{N}(\mathbf{f}_{\setminus i}|f_i)$, which concludes the proof. \square

F.1 Details of Eqn. (6)

$$\begin{aligned}
\text{KL}(\tilde{q}(\mathbf{f})\|\mathcal{N}(\mathbf{f})) &= \int \tilde{q}(\mathbf{f}) \log \frac{\tilde{q}(\mathbf{f}_{\setminus i}|f_i)\tilde{q}(f_i)}{\mathcal{N}(\mathbf{f}_{\setminus i}|f_i)\mathcal{N}(f_i)} d\mathbf{f} \\
&= \int \tilde{q}(f_i) \log \frac{\tilde{q}(f_i)}{\mathcal{N}(f_i)} df_i + \int \tilde{q}(f_i) \int \tilde{q}(\mathbf{f}_{\setminus i}|f_i) \log \frac{\tilde{q}(\mathbf{f}_{\setminus i}|f_i)}{\mathcal{N}(\mathbf{f}_{\setminus i}|f_i)} d\mathbf{f}_{\setminus i} df_i \\
&= \text{KL}(\tilde{q}(f_i)\|\mathcal{N}(f_i)) + \mathbb{E}_{\tilde{q}(f_i)} \left[\text{KL}(\tilde{q}(\mathbf{f}_{\setminus i}|f_i)\|\mathcal{N}(\mathbf{f}_{\setminus i}|f_i)) \right] \\
\tilde{q}(\mathbf{f}_{\setminus i}|f_i) &= \frac{\tilde{q}(\mathbf{f})}{\tilde{q}(f_i)} \propto \frac{p(\mathbf{f})p(y_i|\mathbf{f}_i)\prod_{j \neq i} t_j(\mathbf{f})}{q^{\setminus i}(f_i)p(y_i|f_i)} \\
&= q^{\setminus i}(\mathbf{f}_{\setminus i}|f_i). \tag{23}
\end{aligned}$$

F.2 Details of Eqn. (10)

$$\begin{aligned}
\mathbf{W}_2^2(\tilde{q}(\mathbf{f}), \mathcal{N}(\mathbf{f})) &\equiv \inf_{\pi \in U(\tilde{q}, \mathcal{N})} \mathbb{E}_{\pi} (\|\mathbf{f} - \mathbf{f}'\|_2^2) \\
&= \inf_{\pi \in U(\tilde{q}, \mathcal{N})} \mathbb{E}_{\pi} (\|f_i - f'_i\|_2^2) + \mathbb{E}_{\pi} (\|\mathbf{f}_{\setminus i} - \mathbf{f}'_{\setminus i}\|_2^2) \\
&\stackrel{(a)}{=} \inf_{\pi \in U(\tilde{q}, \mathcal{N})} \mathbb{E}_{\pi_i} [\|f_i - f'_i\|_2^2] + \mathbb{E}_{\pi_{\setminus i|i}} (\|\mathbf{f}_{\setminus i} - \mathbf{f}'_{\setminus i}\|_2^2) \\
&\stackrel{(b)}{=} \inf_{\pi_i} \mathbb{E}_{\pi_i} [\|f_i - f'_i\|_2^2] + \inf_{\pi_{\setminus i|i}} \mathbb{E}_{\pi_{\setminus i|i}} (\|\mathbf{f}_{\setminus i} - \mathbf{f}'_{\setminus i}\|_2^2) \\
&= \inf_{\pi_i} \mathbb{E}_{\pi_i} [\|f_i - f'_i\|_2^2 + \mathbf{W}_2^2(\tilde{q}_{\setminus i|i}, \mathcal{N}_{\setminus i|i})] \\
&\stackrel{(c)}{=} \inf_{\pi_i} \mathbb{E}_{\pi_i} [\|f_i - f'_i\|_2^2 + \mathbf{W}_2^2(q^{\setminus i}_{\setminus i|i}, \mathcal{N}_{\setminus i|i})],
\end{aligned}$$

where the superscript prime indicates that the variable is from the Gaussian \mathcal{N} . In (a), $\pi_i = \pi(f_i, f'_i)$ and $\pi_{\setminus i|i} = \pi(\mathbf{f}_{\setminus i}, \mathbf{f}'_{\setminus i}|f_i, f'_i)$. In (b), the first and the second inf are over $U(\tilde{q}_i, \mathcal{N}_i)$ and $U(\tilde{q}_{\setminus i|i}, \mathcal{N}_{\setminus i|i})$ respectively. (c) is due to $\tilde{q}(\mathbf{f}_{\setminus i}|f_i)$ being equal to $q^{\setminus i}(\mathbf{f}_{\setminus i}|f_i)$ (refer to Eqn. (23)).

F.3 Details of Eqn. (19)

$$\begin{aligned}
&\min_{\mathbf{m}'_{\setminus i}} \text{Eqn. (17)} \\
&= \inf_{\pi_i} \mathbb{E}_{\pi_i} \left[\|f_i - f'_i\|_2^2 + \|\mathbf{a}(f_i - \mu_{\tilde{q}_i}) - \mathbf{a}'(f'_i - m'_i)\|_2^2 \right] \\
&= \inf_{\pi_i} \mathbb{E}_{\pi_i} \left[\|f_i - f'_i\|_2^2 \right] + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 + \|\mathbf{a}'\|_2^2 S'_i - 2\mathbf{a}^\top \mathbf{a}' \mathbb{E}_{\pi_i} (f_i f'_i - \mu_{\tilde{q}_i} m'_i) \\
&= \inf_{\pi_i} \mathbb{E}_{\pi_i} \left[\|f_i - f'_i\|_2^2 \right] + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 + \|\mathbf{a}'\|_2^2 S'_i + \mathbf{a}^\top \mathbf{a}' \mathbb{E}_{\pi_i} \left(\|f_i - f'_i\|_2^2 - f_i^2 - (f'_i)^2 + 2\mu_{\tilde{q}_i} m'_i \right) \\
&= \inf_{\pi_i} \mathbb{E}_{\pi_i} \left[\|f_i - f'_i\|_2^2 \right] + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 + \|\mathbf{a}'\|_2^2 S'_i + \mathbf{a}^\top \mathbf{a}' \mathbb{E}_{\pi_i} \left(\|f_i - f'_i\|_2^2 - (f_i - \mu_{\tilde{q}_i})^2 - \right. \\
&\quad \left. 2f_i \mu_{\tilde{q}_i} + \mu_{\tilde{q}_i}^2 - (f'_i - m'_i)^2 - 2f'_i m'_i + (m'_i)^2 + 2\mu_{\tilde{q}_i} m'_i \right) \\
&= (1 + \mathbf{a}^\top \mathbf{a}') \mathbf{W}_2^2(\tilde{q}_i, \mathcal{N}_i) + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 + \|\mathbf{a}'\|_2^2 S'_i - \mathbf{a}^\top \mathbf{a}' \left(\sigma_{\tilde{q}_i}^2 + \mu_{\tilde{q}_i}^2 + S'_i + (m'_i)^2 - 2\mu_{\tilde{q}_i} m'_i \right) \\
&= (1 + \mathbf{a}^\top \mathbf{a}') \mathbf{W}_2^2(\tilde{q}_i, \mathcal{N}_i) + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 + \|\mathbf{a}'\|_2^2 S'_i - \mathbf{a}^\top \mathbf{a}' \left[\sigma_{\tilde{q}_i}^2 + S'_i + (\mu_{\tilde{q}_i} - m'_i)^2 \right]
\end{aligned}$$

F.4 Details of Eqn. (19)

We first use Proposition 1 to reformulate the L_2 WD $W_2^2(\tilde{q}_i, \mathcal{N}_i)$ as:

$$\begin{aligned}
W_2^2(\tilde{q}_i, \mathcal{N}_i) &= \int_0^1 (F_{\tilde{q}_i}^{-1}(y) - m'_i - \sqrt{2S'_i} \text{erf}^{-1}(2y - 1))^2 dy, \\
&= \int_0^1 (F_{\tilde{q}_i}^{-1}(y) - m'_i)^2 + 2S'_i \text{erf}^{-1}(2y - 1)^2 - 2\sqrt{2S'_i} \text{erf}^{-1}(2y - 1)(F_{\tilde{q}_i}^{-1}(y) - m'_i) dy, \\
&= \int_0^1 (F_{\tilde{q}_i}^{-1}(y) - \mu_{\tilde{q}_i} + \mu_{\tilde{q}_i} - m'_i)^2 dy + S'_i - 2\sqrt{2S'_i} c_{\tilde{q}_i}, \\
&= \sigma_{\tilde{q}_i}^2 + (\mu_{\tilde{q}_i} - m'_i)^2 + S'_i - 2c_{\tilde{q}_i} \sqrt{2S'_i},
\end{aligned}$$

where $F_{\tilde{q}_i}^{-1}(y)$ is the quantile function of $\tilde{q}(f_i)$ and $c_{\tilde{q}_i} \equiv \int_0^1 F_{\tilde{q}_i}^{-1}(y) \text{erf}^{-1}(2y - 1) dy$. Next, we plug this reformulation into Eqn. (19):

$$\begin{aligned}
\text{Eqn. (19)} &= W_2^2(\tilde{q}_i, \mathcal{N}_i) + \mathbf{a}^\top \mathbf{a}' W_2^2(\tilde{q}_i, \mathcal{N}_i) + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 + \|\mathbf{a}'\|_2^2 S'_i - \mathbf{a}^\top \mathbf{a}' \left[\sigma_{\tilde{q}_i}^2 + S'_i + (\mu_{\tilde{q}_i} - m'_i)^2 \right] \\
&= W_2^2(\tilde{q}_i, \mathcal{N}_i) + \mathbf{a}^\top \mathbf{a}' \left[\sigma_{\tilde{q}_i}^2 + (\mu_{\tilde{q}_i} - m'_i)^2 + S'_i - 2c_{\tilde{q}_i} \sqrt{2S'_i} \right] + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 + \|\mathbf{a}'\|_2^2 S'_i \\
&\quad - \mathbf{a}^\top \mathbf{a}' \left[\sigma_{\tilde{q}_i}^2 + S'_i + (\mu_{\tilde{q}_i} - m'_i)^2 \right] \\
&= W_2^2(\tilde{q}_i, \mathcal{N}_i) - 2c_{\tilde{q}_i} \sqrt{2S'_i} \mathbf{a}^\top \mathbf{a}' + \|\mathbf{a}\|_2^2 \sigma_{\tilde{q}_i}^2 + \|\mathbf{a}'\|_2^2 S'_i
\end{aligned}$$

G More Details of EP

We use the expressions $\tilde{q}(\mathbf{f}) = q^{\setminus i}(\mathbf{f})p(y_i|f_i)/Z_{\tilde{q}}$ and $q^{\setminus i}(\mathbf{f}) = q(\mathbf{f})/(t_i(f_i)Z_{q^{\setminus i}})$, and the derivation of $\text{KL}(\tilde{q}(\mathbf{f})\|q(\mathbf{f})) = \text{KL}(\tilde{q}(f_i)\|q(f_i))$ is shown as below:

$$\begin{aligned}
\text{KL}(\tilde{q}(\mathbf{f})\|q(\mathbf{f})) &= \int \tilde{q}(\mathbf{f}) \log \frac{q^{\setminus i}(\mathbf{f})p(y_i|f_i)}{Z_{\tilde{q}}q(\mathbf{f})} d\mathbf{f} \\
&= \int \tilde{q}(\mathbf{f}) \log \frac{q^{\setminus i}(\mathbf{f})p(y_i|f_i)}{Z_{q^{\setminus i}}Z_{\tilde{q}}q^{\setminus i}(\mathbf{f})t_i(f_i)} d\mathbf{f} \\
&= \int \tilde{q}(f_i) \log \frac{p(y_i|f_i)}{Z_{q^{\setminus i}}Z_{\tilde{q}}t_i(f_i)} df_i \\
&= \int \tilde{q}(f_i) \log \frac{q^{\setminus i}(f_i)p(y_i|f_i)}{Z_{q^{\setminus i}}Z_{\tilde{q}}q^{\setminus i}(f_i)t_i(f_i)} df_i \\
&= \int \tilde{q}(f_i) \log \frac{\tilde{q}(f_i)}{q(f_i)} df_i \\
&= \text{KL}(\tilde{q}(f_i)\|q(f_i))
\end{aligned}$$

H Predictive Distributions of Poisson Regression

Given the approximate predictive distribution $f(\mathbf{x}_*) = \mathcal{N}(\mu_*, \sigma_*^2)$ and the relation $g(f) = f^2$, it is straightforward to derive the corresponding $g(\mathbf{x}_*) \sim \text{Gamma}(k_*, c_*)^2$ where the shape k_* and the scale c_* are expressed as [56, 61]:

$$k_* = \frac{(\mu_*^2 + \sigma_*^2)^2}{2\sigma_*^2(2\mu_*^2 + \sigma_*^2)}, \quad c_* = \frac{2\sigma_*^2(2\mu_*^2 + \sigma_*^2)}{\mu_*^2 + \sigma_*^2}.$$

² $\text{Gamma}(x|k, c) = \frac{1}{\Gamma(k)c^k} x^{k-1} e^{-x/c}$.

Furthermore, the predictive distribution of the count value $y \in \mathbb{N}$ can also be derived straightforwardly:

$$\begin{aligned} p(y) &= \int_0^\infty p(g_*)p(y|g_*) dg_* \\ &= \int \text{Gamma}(g_*|k_*, c_*)\text{Poisson}(y|g_*) dg_* \\ &= \frac{c_*^y (c_* + 1)^{-k_* - y} \Gamma(k_* + y)}{y! \Gamma(k_*)} = \text{NB}(y|k_*, c_*/(1 + c_*)), \end{aligned}$$

where $g_* = g(\mathbf{x}_*)$ and NB denotes the negative binomial distribution. The mode is obtained as $\lfloor c_*(k_* - 1) \rfloor$ if $k_* > 1$ else 0.

I Proof of Corollary 2.2

Since the site approximations of both EP and QP are Gaussian, we may analyse the predictive variances using results from the regression with Gaussian likelihood function case, namely the well known Equation (3.61) in [47]:

$$\sigma^2(f_*) = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^\top (K + \tilde{\Sigma})^{-1} \mathbf{k}_*, \quad (24)$$

where $f_* = f(\mathbf{x}_*)$ is the evaluation of the latent function at \mathbf{x}_* and $\mathbf{k}_* = [k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_N)]^\top$ is the covariance vector between the test data \mathbf{x}_* and the training data $\{\mathbf{x}_i\}_{i=1}^N$, K is the prior covariance matrix and $\tilde{\Sigma}$ is the diagonal matrix with elements of site variances $\tilde{\sigma}_i^2$.

After updating the parameters of a site function $t_i(f_i)$, the term $(K + \tilde{\Sigma})^{-1}$ is updated to $(K + \tilde{\Sigma} + (\tilde{\sigma}_{i,\text{new}}^2 - \tilde{\sigma}_i^2) \mathbf{e}_i \mathbf{e}_i^\top)^{-1}$ where $\tilde{\sigma}_{i,\text{new}}$ is the site variance estimated by EP or QP and \mathbf{e}_i is a unit vector in direction i . Using the Woodbury, Sherman & Morrison formula [47, A.9], we rewrite $(K + \tilde{\Sigma} + (\tilde{\sigma}_{i,\text{new}}^2 - \tilde{\sigma}_i^2) \mathbf{e}_i \mathbf{e}_i^\top)^{-1}$ as

$$\begin{aligned} &(K + \tilde{\Sigma} + (\tilde{\sigma}_{i,\text{new}}^2 - \tilde{\sigma}_i^2) \mathbf{e}_i \mathbf{e}_i^\top)^{-1} \\ &\equiv (A^{-1} + (\tilde{\sigma}_{i,\text{new}}^2 - \tilde{\sigma}_i^2) \mathbf{e}_i \mathbf{e}_i^\top)^{-1} \\ &= A - A \mathbf{e}_i [(\tilde{\sigma}_{i,\text{new}}^2 - \tilde{\sigma}_i^2)^{-1} + \mathbf{e}_i^\top A \mathbf{e}_i]^{-1} \mathbf{e}_i^\top A \\ &\equiv A - \mathbf{s}_i [(\tilde{\sigma}_{i,\text{new}}^2 - \tilde{\sigma}_i^2)^{-1} + A_{ii}]^{-1} \mathbf{s}_i^\top \\ &= A - \frac{1}{(\tilde{\sigma}_{i,\text{new}}^2 - \tilde{\sigma}_i^2)^{-1} + A_{ii}} \mathbf{s}_i \mathbf{s}_i^\top \end{aligned}$$

where $A = (K + \tilde{\Sigma})^{-1}$ and \mathbf{s}_i is the i 'th column of A . Putting the above expression into Equation (24), we have that the predictive variance is updated according to:

$$\sigma_{\text{new}}^2(f_*) = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^\top A \mathbf{k}_* + \frac{1}{(\tilde{\sigma}_{i,\text{new}}^2 - \tilde{\sigma}_i^2)^{-1} + A_{ii}} \mathbf{k}_*^\top \mathbf{s}_i \mathbf{s}_i^\top \mathbf{k}_*.$$

In EP and QP, the first two terms on the r.h.s. of the above equation are equivalent. As the site variance provided by QP is less or equal to that by EP, *i.e.*, $\tilde{\sigma}_{i,\text{QP}}^2 \leq \tilde{\sigma}_{i,\text{EP}}^2$, the third term on the r.h.s. for QP is less or equal to that for EP. Therefore, the predictive variance of QP is less or equal to that of EP: $\sigma_{\text{QP}}^2(f_*) \leq \sigma_{\text{EP}}^2(f_*)$.

J Lookup Tables

To speed up updating variances σ_{QP}^2 in QP, we pre-compute the integration in Equation (5) over a grid of cavity parameters μ and σ , and store the results into lookup tables. Consequently, each update step obtains σ_{QP}^2 simply based on the lookup tables. Concretely, for the GP binary classification, we compute Equation (5) with μ , σ and y varying from -10 to 10, 0.1 to 10 and $\{-1, 1\}$ respectively. μ and σ vary in a linear scale and a log10 scale respectively, and both have a step size of 0.001. The resulting lookup tables has a size of 20001×2001 . In a similar way, we make the lookup table

Algorithm 1 Expectation (Quantile) Propagation

Input: $p(\mathbf{f}), p(y_i|f_i), t_i(f_i), i = 1, \dots, N, \theta$ **Output:** $q(\mathbf{f})$ approximate posterior
1: **repeat**
2: compute $q(\mathbf{f}) \propto p(\mathbf{f}) \prod_i t_i(f_i)$ by (1)
3: **repeat**
4: **for** $i = 1$ to N **do**
5: compute $q^{i}(f_i) \propto q(f_i)/t_i(f_i)$ cavity
6: compute $\tilde{q}(f_i) \propto q^{i}(f_i)p(y_i|f_i)$ tilted
7: **if** EP **then**
8: $t_i(f_i) \propto \text{proj}_{\text{KL}}[\tilde{q}(f_i)]/q^{i}(f_i)$ by (3)(4)
9: **else if** QP **then**
10: $t_i(f_i) \propto \text{proj}_{\text{W}}[\tilde{q}(f_i)]/q^{i}(f_i)$ by (5)(4)
11: **end if**
12: update $q(\mathbf{f}) \propto p(\mathbf{f}) \prod_i t_i(f_i)$ by (1)
13: **end for**
14: **until** convergence
15: $\theta = \text{argmax}_{\theta} \log q(\mathcal{D})$ by (2)
16: **until** convergence
17: **return** $q(\mathbf{f})$

for the Poisson regression. In the experiments, we exploit the linear interpolation to fit σ_{QP}^2 given $\mu \in [-10, 10]$ and $\sigma \in [0.1, 10]$, and if μ and σ lie out of the lookup table, σ_{QP}^2 is approximately computed by the EP update formula, i.e., $\sigma_{\text{QP}}^2 \approx \sigma_{\text{EP}}^2$. On Intel(R) Xeon(R) CPU E5-2680 v4 @ 2.40GHz, we observe the running time of EP and QP is almost the same.