

Bias & variance:

$$\text{Bias of } \hat{s} : b(s) = E[\hat{s} | s] - s$$

(defined for any parameter θ of a statistical model)

$$\text{Var}(\hat{s} | s) = E[(\hat{s} - E[\hat{s} | s])^2 | s]$$

$$E[(\hat{s} - s)^2] = E[(\hat{s} - E[\hat{s} | s] + b(s))^2 | s]$$

$$= \underbrace{\text{Var}(\hat{s} | s)}_{\text{variance}} + \underbrace{b^2(s)}_{\text{bias}}$$

$b(s) = 0$: unbiased estimator

Fisher information (scalar case):

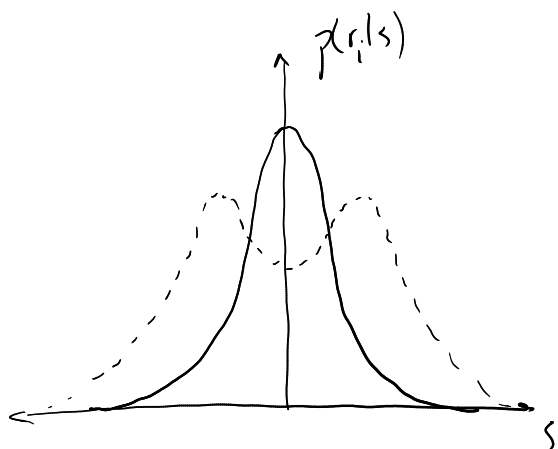
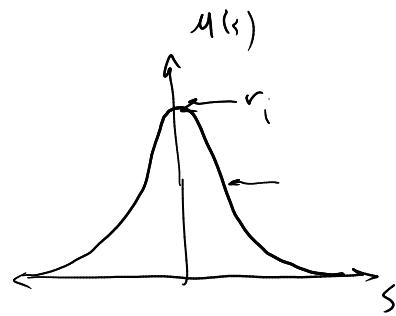
$$I_F(s) = - E \left[\frac{\partial^2 \log p(r|s)}{\partial s^2} \right]_{p(r|s)} \quad (1)$$

$$= - \int dr p(r|s) \frac{\partial^2 \log p(r|s)}{\partial s^2} \quad \left[\begin{array}{l} \text{or} \\ E \left[\left(\frac{\partial}{\partial s} \log p(r|s) \right)^2 \right]_{p(r|s)} \end{array} \right] \quad (2)$$

Example:

$$\lambda_i = \mu(s) \quad \uparrow \text{ tuning curve}$$

$$r_i \sim \text{Pois}(\lambda_i)$$



$$p(r_i | s) = \frac{\mu(s)^{r_i} e^{-\mu(s)}}{r_i!}$$

$$\log p(r_i | s) = r_i \log \mu(s) - \mu(s) - \log r_i!$$

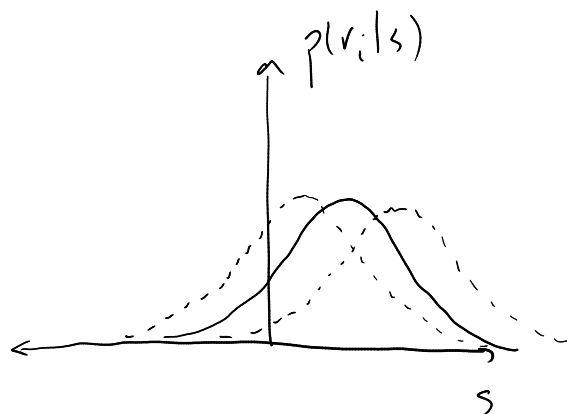
$$\frac{\partial}{\partial s} \log p(r_i | s) = r_i \frac{\mu'(s)}{\mu(s)} - \mu'(s) = \mu'(s) \left[\frac{r_i}{\mu(s)} - 1 \right]$$

$$E[r_i^2] = (\mu'(s))^2 E \left[\frac{r_i^2}{\mu(s)} - \frac{2r_i}{\mu(s)} + 1 \right]$$

$$E[r_i^2] = \mu(s)(1 + \mu(s))$$

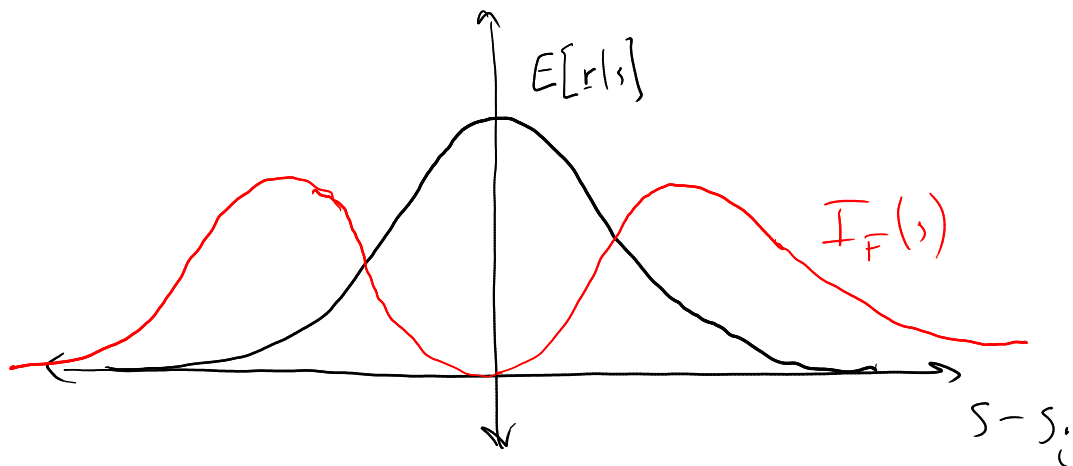
$$E[r_i] = \mu(s)$$

$$= \frac{(\mu'(s))^2}{\mu(s)} \left(\frac{1}{\mu(s)} + 1 - 2 + 1 \right)$$



$$\lambda_i = r_{\max} \exp\left(\frac{-(s_i - s)^2}{2\sigma_i^2}\right) \quad r_i \sim \text{Pois}(\lambda_i)$$

$$I_F(s) = \frac{r_{\max} (s_i - s)^2}{\sigma_i^4} \exp\left(\frac{-(s_i - s)^2}{2\sigma_i^2}\right)$$



Intuition (using def 1): Expected curvature of log-likelihood fn.

Intuition (using def. 2): "Score": $\frac{d}{ds} \log p(r|s)$

How much does log-likelihood of observing r change when s is varied?

$$\begin{aligned} E\left[\frac{d}{ds} \log p(r|s)\right]_{p(r|s)} &= \int dr \cancel{p(r|s)} \cdot \frac{\frac{d}{ds} p(r|s)}{\cancel{p(r|s)}} \\ &= \frac{d}{ds} \int dr p(r|s) = \frac{d}{ds} (1) = 0. \end{aligned}$$

$$\text{So } \text{Var}(\text{score}) = E[\text{score}^2] = I_F(s)$$

Properties:

1) Local (dependent on value of s)

2) If r_i, r_j independent, $I_F(s) = I_F^i(s) + I_F^j(s)$

$$\log p(\underline{r}|s) = \log p(r_i|s) + \log p(r_j|s)$$

3) More generally, $I_F^{X,Y} = I_F^X + I_F^{Y|X}$

5) Related

to variance of unbiased estimator

(Cramér-Rao bound, later).

4) Dependent on stimulus parameterization:

$$\text{If } u = f(s), I_F(s) = I_F(u) \left(\frac{du}{ds} \right)^2$$

Equivalence of 2 definitions:

$$\frac{d^2}{ds^2} \log p(\underline{r}|s) = \frac{d}{ds} \left[\frac{d}{ds} \log p(\underline{r}|s) \right]$$

$$= \frac{d}{ds} \left[\frac{1}{p(\underline{r}|s)} \frac{d}{ds} p(\underline{r}|s) \right]$$

$$= \frac{1}{p(\underline{r}|s)} \frac{d^2}{ds^2} p(\underline{r}|s) - \left(\frac{\frac{d}{ds} p(\underline{r}|s)}{p(\underline{r}|s)} \right)^2$$

$$E[\cdot] = \int d\underline{r} \frac{d^2}{ds^2} p(\underline{r}|s)$$

$$= \frac{d^2}{ds^2} \int d\underline{r} p(\underline{r}|s) = 0$$

$$= \left[\frac{d}{ds} \log p(\underline{r}|s) \right]^2$$

$$\text{So } I_F(s) = E \left[\left(\frac{d}{ds} \log p(\underline{r}|s) \right)^2 \right]_{p(\underline{r}|s)} \quad (2)$$

Cramér-Rao bound: If $\hat{\theta}$ is an unbiased estimator of a parameter θ of a distribution $p(r|\theta)$, then

$$\text{Var}(\hat{\theta}|\theta) \geq \frac{1}{I_F(\theta)}$$

Recall unbiased $\Rightarrow E[\hat{\theta}|\theta] - \theta = 0$.

Proof: Score $A = \frac{\partial}{\partial \theta} \log p(r|\theta)$. From last time, $E[A] = 0$.

Also let $B = \hat{\theta} - \theta$ $E[B] = 0$

Use Cauchy-Schwartz inequality:

$$E[AB]^2 \leq E[A^2]E[B^2]$$

$$E[A^2] = I_F(\theta)$$

$$E[B^2] = \text{Var}(\hat{\theta}|\theta)$$

$$\begin{aligned} E[AB] &= \int dr \left[(\hat{\theta} - \theta) \frac{\partial}{\partial \theta} \log p(r|\theta) \right] p(r|\theta) \\ &= \underbrace{\int dr \hat{\theta} \frac{\partial}{\partial \theta} \log p(r|\theta) p(r|\theta)}_{(1)} - \underbrace{\int dr \theta \frac{\partial}{\partial \theta} \log p(r|\theta) p(r|\theta)}_{(2)} \end{aligned}$$

$$(1) = \frac{d}{d\theta} \int d\underline{r} \hat{\theta} p(\underline{r}|\theta) = \frac{d}{d\theta} (\theta) = 1 \quad (\text{note } \hat{\theta} \text{ is a function of } \underline{r}, \text{ not } \theta \text{ explicitly})$$

$$(2) = \theta \int d\underline{r} \frac{d}{d\theta} \log p(\underline{r}|\theta) p(\underline{r}|\theta) = \theta E\left[\frac{d}{d\theta} \log p(\underline{r}|\theta)\right] = \theta E[\text{Score}] = 0$$

$$\Rightarrow I \leq \bar{I}_F(\theta) \times \text{Var}(\hat{\theta}|\theta)$$

$\hookrightarrow = 0$ from last time.

$$\Rightarrow \text{Var}(\hat{\theta}|\theta) \geq 1/\bar{I}_F(\theta)$$

Generalizations:

1) Biased estimators:

If $\hat{\theta}$ has bias $b(\hat{\theta}|\theta)$, then

$$\text{Var}(\hat{\theta}|\theta) \geq \frac{[1 + b'(\hat{\theta}|\theta)]^2}{\bar{I}_F(\theta)}$$

$$b' = \frac{d}{d\theta} b(\hat{\theta}|\theta)$$

2) Multivariate case

Fisher information matrix:

$$[\bar{I}_F(\underline{s})]_{ij} = E\left[\left(\frac{d}{ds_i} \log p(\underline{r}|\underline{s})\right) \left(\frac{d}{ds_j} \log p(\underline{r}|\underline{s})\right) \mid \underline{s}\right]$$

$$= -E\left[\frac{d}{ds_i ds_j} \log p(\underline{r}|\underline{s}) \mid \underline{s}\right]$$

$$\underline{a}^T \text{Cov}(\hat{\underline{S}} | \underline{s}) \underline{a} \geq \underline{a}^T \left[\mathbf{I}_F(\underline{s}) \right]^{-1} \underline{a} \quad \forall \underline{a}$$

(diagonal elements \Rightarrow univariate case)

Covariance matrices

$$[\Sigma(\underline{x})]_{ij} = \text{Cov}(x_i, x_j) = E \left[(x_i - E[x_i]) (x_j - E[x_j]) \right]$$

$$\Sigma = \begin{pmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \dots \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) & \dots \\ \vdots & \vdots & \ddots \\ & & & \text{Var}(x_n) \end{pmatrix}$$

Properties:

1) Symmetric

2) Positive-semidefinite: $\underline{a}^T \Sigma \underline{a} \geq 0$

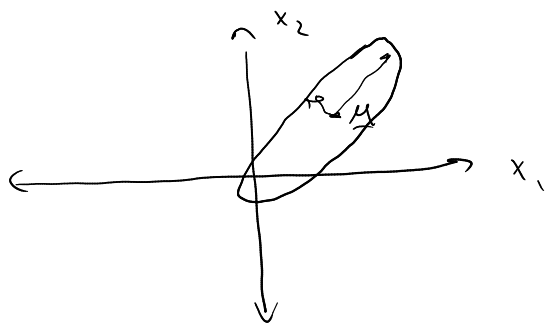
3) E-evals, e-vecs related to PCA. Note

Σ, Σ^{-1} have same e-vecs (see with SVD).

Multivariate normal distribution:

$$E[\underline{x}] = \underline{\mu}, \quad \text{Cov}(\underline{x}) = \Sigma \quad \begin{array}{l} \underline{x} \sim \mathcal{N}(\underline{\mu}, \Sigma) \\ \underline{x} \in \mathbb{R}^N \end{array}$$

$$p(\underline{x}) = \frac{1}{\sqrt{(2\pi)^N \det(\Sigma)}} \exp\left(-\frac{1}{2} (\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu})\right)$$



Σ diagonal: independent.

Linear Fisher information (Moreno-Bote et al., Nat Neurosci 2014)

Suppose N neurons w/ tuning curves $E[r_i | s] = \mu_i(s)$,

$$i = 1 \dots N. \text{ Let } \underline{\mu}(s) = \begin{bmatrix} \mu_1(s) \\ \mu_2(s) \\ \vdots \\ \mu_N(s) \end{bmatrix}$$

$$\text{and } \underline{\mu}'(s) = \frac{\partial}{\partial s} \underline{\mu}(s). \text{ (tuning curve slope).}$$

Also assume that $r_i(s) \sim \mathcal{N}(\underline{\mu}(s), \Sigma(s))$.

Note: $\Sigma(s)$ noise covariance matrix. Correlation:

$$\rho_{ij} = \frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii} \Sigma_{jj}}}$$

Can find: $I_F(\cdot) = (\underline{\mu}')^T \Sigma^{-1} \underline{\mu}' +$

$$\frac{1}{2} \text{Tr} \left[\left(\Sigma' \Sigma^{-1} \right)^2 \right]$$

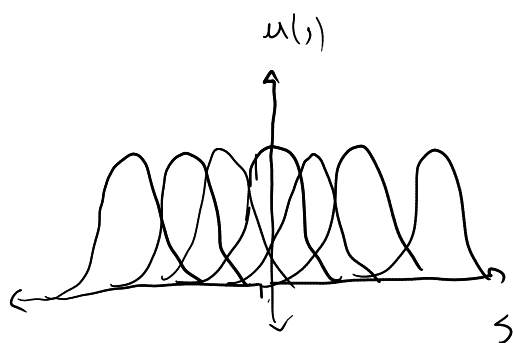
$\Sigma' = \frac{d}{ds} \Sigma(s)$: Stimulus-dependent noise correlations.

If we consider a ^{"locally optimal"} linear estimator $\hat{s} = \sum_i w_i r_i$,

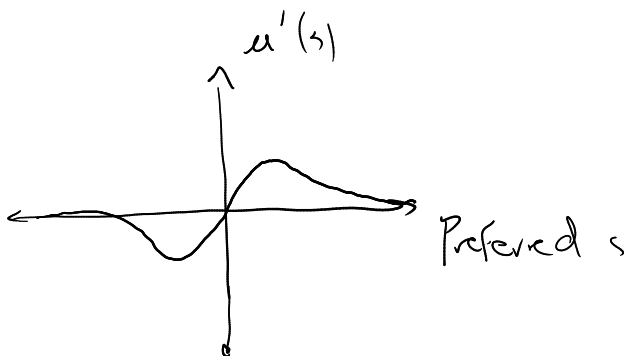
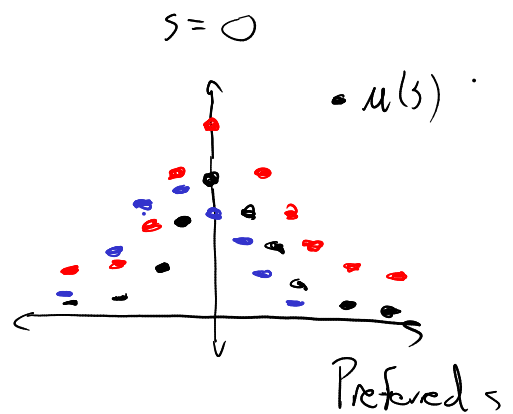
$$\text{Var}(\hat{\theta}|\theta) \geq \frac{1}{I_{\text{LOLE}}(s)}, \text{ where}$$

$$I_{\text{LOLE}}(s) = (\underline{\mu}')^T \Sigma^{-1} \underline{\mu}'$$

Also valid if $p(\underline{r}|s)$ is of exponential family w/ linear sufficient statistics.



\Rightarrow



Consequences:

1) Correlations in the direction of $(\underline{\mu}')^T \underline{\mu}'$ limit information as $N \rightarrow \infty$. "Differential correlations", identical tuning curves \Rightarrow pos. mean correlations.

If $\Sigma = \Sigma_0 + \varepsilon (\underline{\mu}')^T \underline{\mu}'$, where $I_0(\cdot) \rightarrow \infty$ as $N \rightarrow \infty$ if $\Sigma = \Sigma_0$, then

$$I = \frac{I_0}{1 + \varepsilon I_0}. \quad \text{As } N \rightarrow \infty, \quad I \rightarrow \frac{1}{\varepsilon}$$

2) Differential correlations may be hidden. Instead, use decoder.