

Let  $r_i = \underbrace{f_i(x)}_{\text{fuzzy}} + \underbrace{\epsilon_i}_{\text{noise model}}$

$\Rightarrow$  Given  $r_i$ , how much information do we have about  $x$ ?

Translates to: how much linear Fisher info is in the population?

We can think about this in terms of a local decoder  $w$ :

$$\hat{x} = x_0 + w^T (r - f(x_0))$$

We want this estimate to be unbiased:

$$\langle \hat{x} \rangle = x$$

$$\Rightarrow x_0 + w^T (\langle r \rangle - f(x_0)) = x$$

$$x_0 + w^T (f(x) - f(x_0))$$

$$\approx x_0 + w^T (f(x_0) + f'(x_0)(x - x_0) - f(x_0))$$

$$= x_0 + w^T f'(x_0) (x - x_0)$$

~~$$\Rightarrow w^T f'(x_0) (x - x_0) = x - x_0$$~~

$$\Rightarrow \underline{\underline{w^T f'(x_0) = 1}}$$

• Minimum variance:

$$\begin{aligned} \text{Var} \{ \hat{x} \} &= w^T \text{Var} [ r - f(x_0) ] w \\ &= w^T \left( \langle (r - f(x_0))(r - f(x_0))^T \rangle \right. \\ &\quad \left. - \langle r - f(x_0) \rangle \langle r - f(x_0) \rangle^T \right) w \\ &= w^T \left[ \langle r r^T \rangle - \langle r \rangle \langle r \rangle^T \right] w \\ &= w^T \left[ \langle (f(x) + \epsilon)(f(x) + \epsilon)^T \rangle - f(x) f(x)^T \right] w \\ &= w^T \underbrace{\langle \epsilon \epsilon^T \rangle}_R w \\ &= \underline{\underline{w^T R w}} \end{aligned}$$

So, now find  $w$  that

minimizes  $\text{Var} \{ \hat{x} \} = w^T R w$   
subject to  $w^T f'(x_0) = 1$

$$\Rightarrow \frac{\partial}{\partial w} \left[ w^T R w + \lambda (w^T f'(x_0) - 1) \right] = 0$$

$$\Rightarrow \alpha \mathbf{1}^T \mathbf{w} = 1$$

$$\Leftrightarrow \underline{\mathbf{w}} = \alpha \mathbf{R}^{-1} \mathbf{f}'(x_0)$$

$$\Rightarrow \mathbf{w}^T \mathbf{f}'(x_0) = \alpha \mathbf{f}'(x_0)^T \mathbf{R}^{-1} \mathbf{f}'(x_0) = 1$$

$$\Leftrightarrow \alpha = \frac{1}{\mathbf{f}'(x_0)^T \mathbf{R}^{-1} \mathbf{f}'(x_0)}$$

$$\Rightarrow \boxed{\mathbf{w} = \frac{\mathbf{R}^{-1} \mathbf{f}'(x_0)}{\mathbf{f}'(x_0)^T \mathbf{R}^{-1} \mathbf{f}'(x_0)}}$$

Giving us:

$$\text{Var}[\hat{x}] = \mathbf{w}^T \mathbf{R} \mathbf{w}$$

$$= \frac{1}{\mathbf{f}'(x_0)^T \mathbf{R}^{-1} \mathbf{f}'(x_0)}$$

We then define the Fisher Information as

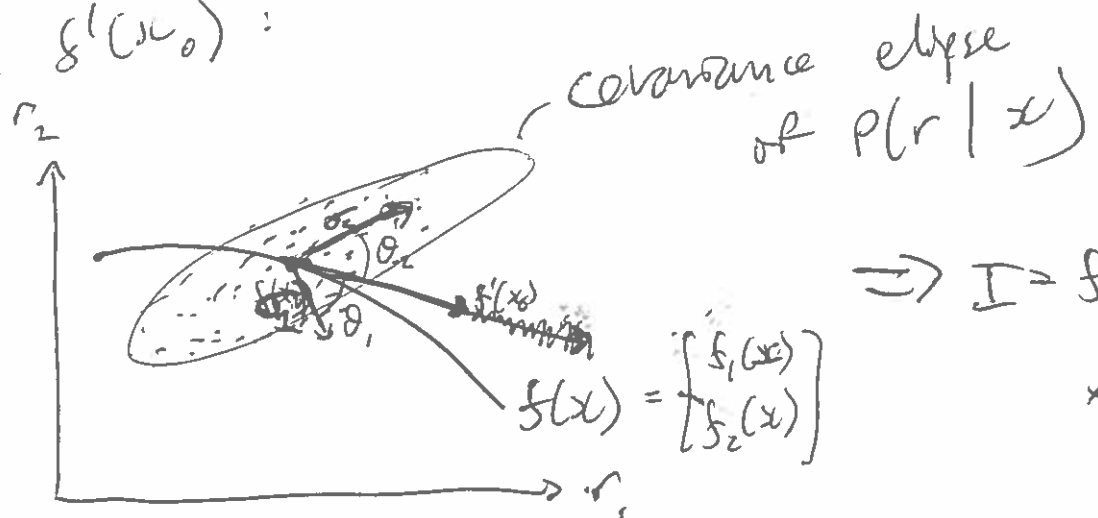
$$\underline{I} = \frac{1}{\text{Var}[\hat{x}]} = \mathbf{f}'(x_0)^T \mathbf{R}^{-1} \mathbf{f}'(x_0)$$

~~The fact, this is dependent on the angle~~

Geometrically, the decoder is computing the gradient at  $x_0$  and adjusting it by the direction of variance in the distribution over  $r$  (see slides).

Intuitively, then, the amount of information ~~should~~ in the code should depend on those directions of highest variance, ~~and variance along  $\mathbf{f}'(x_0)$  can't be averaged away~~ ~~so should decrease  $\mathbf{f}'(x_0)$~~ . Also picture it in your head: variance perpendicular to  $\mathbf{f}'(x_0)$  shouldn't hurt  $\hat{x}$ .

Indeed, it turns out that we can rewrite the linear Fisher Info in terms of angles & w directions of maximal variance (e-vectors of  $R$ ) and  $f'(x_0)$ :



$$f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \end{bmatrix}$$

$$\Rightarrow I = f'(x_0)^T f'(x_0) \times \left( \frac{\cos^2 \theta_1}{\sigma_1^2} + \frac{\cos^2 \theta_2}{\sigma_2^2} \right)$$

Pf. Recall that any symmetric matrix can be rewritten as:

$$A v_k = \lambda_k v_k, \quad v_k^T v_j = \delta_{jk} \quad \left( \sum_k \lambda_k v_k v_k^T \right) v_j = \lambda_j v_j$$

$$\Rightarrow A = \sum_k \lambda_k v_k v_k^T$$

where  $\lambda_k, v_k$  are the e-values, vectors of  $A$

~~Thus, we can write~~  
 ~~$A = \sum_k \lambda_k v_k v_k^T$~~   
 ~~$A^{-1} = \sum_k \frac{1}{\lambda_k} v_k v_k^T$~~

Then,

$$\left( \sum_k \lambda_k^{-1} v_k v_k^T \right) A$$

$$= \sum_{j,k} \frac{1}{\lambda_k} v_k v_k^T \lambda_j v_j v_j^T$$

$$= \sum_k v_k v_k^T = I$$

$$\Rightarrow A^{-1} = \sum_k \frac{1}{\lambda_k} v_k v_k^T$$

(since  $\sum_k v_k v_k^T = I$ )  
 $= [v_1 \dots v_n] \begin{bmatrix} \lambda_1^{-1} & & \\ & \dots & \\ & & \lambda_n^{-1} \end{bmatrix} \begin{bmatrix} -v_1^T \\ \dots \\ -v_n^T \end{bmatrix}$   
 $= I$   
 since  $Q$  is orthogonal

We thus have

$$R^{-1} = \sum_k \frac{1}{\sigma_k} v_k v_k^T$$

$$\Rightarrow I = \sum_k \frac{(f'(x_0)^T v_k)^2}{\sigma_k^2}$$

$$= \sum_k \frac{|f'(x_0)|^2 |v_k|^2 \cos^2 \theta_k}{\sigma_k^2}$$

angle btw  $f'(x_0)$  and  $v_k$   
( $v_k$  is e-vector of  $R^{-1}$ )

$$= \left( \sum_k \frac{\cos^2 \theta_k}{\sigma_k^2} \right) \underbrace{f'(x_0)^T f'(x_0)}_{= |f'(x_0)|^2}$$

We now note some properties:

~~the Jacobian matrix is always big~~

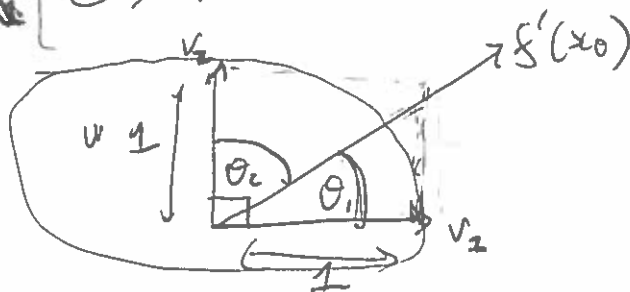
•  $f'(x_0)^T f'(x_0) \in \mathbb{N} \sim \text{big}$

•  $\sum_k \cos^2 \theta_k = 1$

since  $v_1 \perp v_2$ , ~~we have~~  $\Rightarrow$  think about 2-D case:

$$\cos \theta_2 = \cos\left(\frac{\pi}{2} - \theta_1\right) = \sin \theta_1$$

$$\Rightarrow \cos^2 \theta_1 + \cos^2 \theta_2 = \cos^2 \theta_1 + \sin^2 \theta_1 = \underline{\underline{1}}$$



$$\bullet \sum_k \sigma_k^2 \ll N$$

assuming all  $\sigma_k^2$  are of similar size (i.e. spherical noise) \*

$$\Rightarrow \sum_k \frac{\cos^2 \theta_k}{\sigma_k^2} = \text{convex combination of } \sigma_k^2 \text{'s} \sim \mathcal{O}(1)$$

(approximately independent of  $N$ )

$$\Rightarrow I \ll N$$

assuming  $\sigma_k^2$  are independent (spherical noise)

or weakly correlated such that  $\sigma_k^2$ 's still similar size

What about when noise is not independent?

→ some of the  $\sigma_k^2$  then become proportional to  $N$

Then, if  $f'(x_0)$  is parallel to  $v_k$  with  $\sigma_k^2 \ll N$ ,  $\cos^2 \theta_j = \delta_{jk}$  and thus

$I \ll \mathcal{O}(1)$  will be constant w.r.t.  $N$  (noise can't be averaged away)

But does this ever happen!

It seems unlikely that the  $v_k$ 's with largest covariance ( $\sigma_k^2$ ) would line up exactly with  $f'(x_0)$  ...

⇒ Shamir & Sompolinsky 2006

Eden, Berens, Tolvas & Bethge 2011

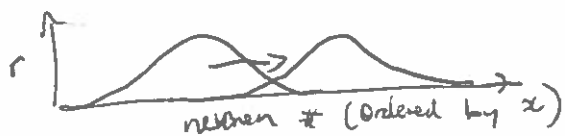
NO

(more specifically, whenever tuning curves are not exactly the same (e.g. amplitude, width, etc.) across neurons, which is the case in the brain!)

So, under realistic conditions (i.e.

heterogeneous tuning curves, lots of neurons) the directions of highest covariance do not exactly line up with the derivative of the tuning curve and thus noise can be averaged away with large N.

In other words, the only correlations that make  $\Sigma$  independent of  $N$  are those that push the responses of all the neurons in the same direction: we call these differential correlations







$$\frac{(m_1 + m_2)}{(1+c)^2 m_1 m_2 + 8}$$

---

$$16(1+c)m_1 m_2$$

## Correlations Study

Shallen et al 1994:

correlations bound SAR  
in large N limit  
→ homogeneous tuning curves

Dayan & Abbott '98:

correlations are good  
for homogeneous tuning with  
equally spaced preferences  
→ except when correlations  
a function of pref stim.

Shannon & Sompolsky 2006

Olsh et al. 2011:

For heterogeneous tuning curves,  
correlations don't matter

Moore - Bote et al 2014:

The only conclusions  
that matters are  
differential conclusions

→ Also, information can't  
be infinite!

(finite input info,  
infinite computation)

→ Therefore, there must  
be differential conclusions

---

(But hard to see,  
generally overlooked  
by <sup>other</sup> reverse conclusions)