

ENCODING & DECODING I

0. Introduction

- Let's first start describing a dataset that ~~and~~ represents a typical dataset we would encounter in neuroscience.



- In this and the upcoming ~~and~~ lecture we are going to see the two most important branches of supervised learning: regression and classification.
- Regression and classification can be defined as predicting a dependent variable y from an independent D dimensional variable x .

- The difference between regression and classification relies on the fact that in regression the dependent variable is continuous whereas in classification is discrete.

In general we are going to assume that:

$$y = f(\vec{x})$$

- Restrict a bit more our assumption and say:

$$y = \vec{w}^T \vec{\phi}(\vec{x}) + w_0$$

where $\phi_i(\vec{x})$ are called the basis functions and represent in general an arbitrary function of x .

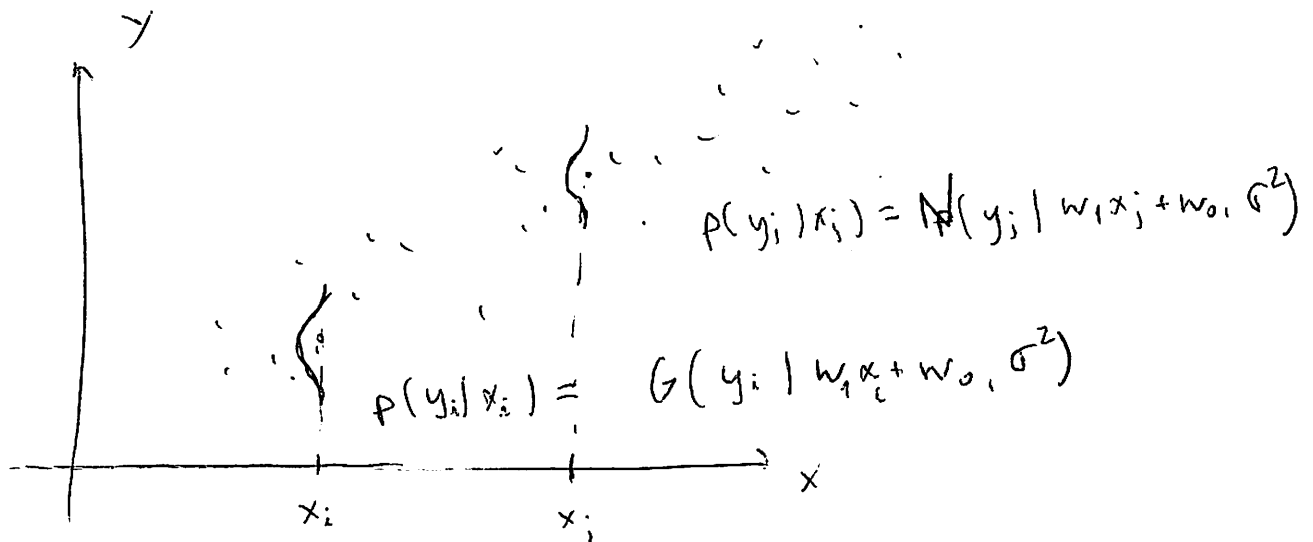
Here we are going to focus only on linear regression and therefore $\vec{\phi}(\vec{x}) = \vec{x}$.

$$y = \vec{w}^T \vec{x} + w_0$$

Let's assume for instance

$$\rightarrow \boxed{r_i = \vec{w}^T \vec{B} + w_0}$$

1. Maximum likelihood and least squares for linear regression:



Our model assumes:

$$y = \vec{w}^T \vec{x} + w_0 + \sigma \tilde{z} = \vec{w}^T \vec{x} + \sigma \tilde{z}$$

where $\vec{x} = (\vec{x}, 1)$

$$\mathbb{E}[y|x] = \int y p(y|x) dy = \vec{w}^T \vec{x}$$

Consider now a set of $\{y_1, \dots, y_N\}, \{\vec{x}_1, \dots, \vec{x}_N\}$

The likelihood becomes:

$$p(\vec{y} | \vec{X}, \vec{w}) = \prod_i p(y_i | \vec{x}_i, \vec{w}) = \prod_i \mathcal{N}(y_i | \vec{w}^T \vec{x}_i, \sigma^2)$$

iid

$$\begin{aligned} \rightarrow \ln p(\vec{y} | \vec{X}, \vec{w}) &= \sum_i \ln p(y_i | \vec{x}_i, \vec{w}) = -\frac{N}{2} \ln(\sigma^2) - \frac{N}{2} \ln(2\pi) \\ &\quad - \frac{1}{2\sigma^2} \sum_i (y_i - \vec{w}^T \vec{x}_i)^2 \end{aligned}$$

$$\Rightarrow \boxed{E(\vec{w}) = \frac{1}{2} \sum_i (y_i - \vec{w}^T \vec{x}_i)^2} \quad \leftarrow \text{Only one \textit{arg} depending on } \vec{w}$$

Least squares cost function

Maximize log prob. under Gaussian noise is equivalent to minimize least squares cost function.

$$\nabla_{\vec{w}} \ln p(\vec{y} | \vec{X}, \vec{w}) = \frac{1}{\sigma^2} \sum_i (y_i - \vec{w}^T \vec{x}_i) \vec{x}_i^T = 0$$

$$\sum_i y_i \vec{x}_i^T - \sum_i \vec{w}^T \vec{x}_i \vec{x}_i^T = 0 \Rightarrow \vec{w}^T \sum_i \vec{x}_i \vec{x}_i^T = \sum_i y_i \vec{x}_i^T$$

$$\boxed{\vec{w}_{ML} = (X^T X)^{-1} X^T \vec{y}}$$

let's have a look to the bias explicitly:

$$\frac{\partial}{\partial w_0} \ln P(\vec{y} | X) \stackrel{\sigma^2}{=} \frac{\partial}{\partial w_0} \frac{1}{\sigma^2} \sum (y_i - (\vec{w}^T \vec{x}_i + w_0))^2 = 0$$

$$\cancel{\sum (y_i - (\vec{w}^T \vec{x}_i + w_0))^2} = \frac{1}{\sigma^2} \sum (y_i - \vec{w}^T \vec{x}_i + w_0) = 0$$

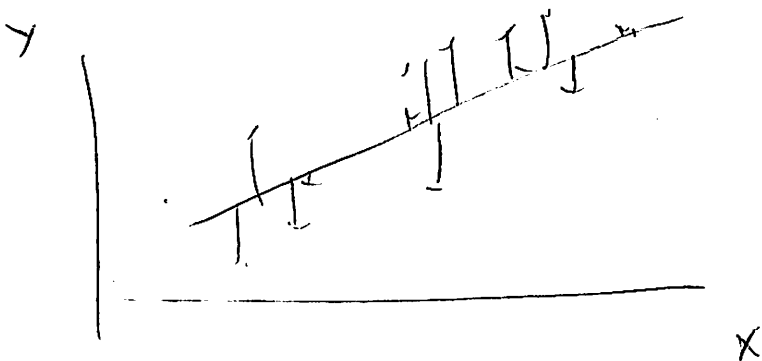
$$\sum y_i - \sum \vec{w}^T \vec{x}_i - N w_0 = 0 \Rightarrow \boxed{w_{0,ML} = \frac{\sum y_i}{N} - \vec{w}_{ML}^T \frac{\sum \vec{x}_i}{N}}$$

$$w_{0,ML} = \bar{y} - \vec{w}_{ML}^T \bar{\vec{x}}$$

Solution for σ^2

$$\sigma^2 = \frac{1}{N} \sum_i (y_i - \vec{w}_{ML}^T \vec{x}_i)^2$$

which is the sum of residuals square



2. Regularized least squares

In the previous section we derived the likelihood function:

$$N(\vec{y} | \vec{m}_0, \sigma^2 \mathbb{I})$$

$$\vec{m}_0 = \vec{0}, \quad \sigma^2 = \frac{1}{\sigma^2} X^T X$$

$$p(\vec{y} | \vec{x}, \vec{w}, \sigma^2) = \prod_i p(y_i | x_i, \vec{w}, \sigma^2) = \prod_i N(y_i | \vec{w}^T x_i, \sigma^2)$$

To make things easier, let's skip the dependence on \vec{x} and σ^2 .

$$p(\vec{y} | \vec{w}) = N(\vec{y} | \vec{w})$$

Bayes Rule: $p(a|b) = \frac{p(b|a)p(a)}{p(b)} \propto p(b|a)p(a)$

$$\Rightarrow \underbrace{p(\vec{w} | \vec{y})}_{\text{Gaussian}} \propto \underbrace{p(\vec{y} | \vec{w})}_{\text{Gaussian}} \underbrace{p(\vec{w})}_{\text{Gaussian}}$$

Gaussian is the conjugate prior of itself

Let's assume $p(\vec{w}) = N(\vec{w} | \vec{m}_0, S_0)$

Then

Exercise

$$p(\vec{w} | \vec{y}) = N(\vec{w} | \vec{m}_N, S_N)$$

$$\text{where } \begin{cases} \vec{m}_N = S_N (S_0^{-1} \vec{m}_0 + \frac{1}{\sigma^2} X^T \vec{y}) \\ S_N^{-1} = S_0^{-1} + \frac{1}{\sigma^2} X^T X \end{cases}$$

Let's assume $S_0 = \alpha \mathbb{I}$, then: $\begin{cases} \vec{m}_N = S_N (\alpha \vec{m}_0 + \frac{1}{\sigma^2} X^T \vec{y}) \\ S_N^{-1} = \alpha \mathbb{I} + \frac{1}{\sigma^2} X^T X \end{cases}$

$$p(\vec{w} | \alpha) = N(\vec{w} | \vec{0})$$

If the prior is infinitely broad, $\alpha \rightarrow 0$, then:

$$\vec{m}_N = (X^T X)^{-1} X^T \vec{y} \Rightarrow \vec{w}_{ML} \quad \& \text{likelihood maximization.}$$

• let's now assume $p(\vec{w} | \alpha) = \mathcal{N}(\vec{w} | \vec{0}, \alpha^{-1} \mathbb{I})$

$$\begin{cases} \vec{m}_N = \frac{1}{\sigma^2} S_N X^T \vec{y} \\ S_N^{-1} = \alpha \mathbb{I} + \frac{1}{\sigma^2} X^T X \end{cases} \quad \hookrightarrow p(\vec{w} | \vec{y}) = \mathcal{N}(\vec{w} | \vec{m}_N, S_N)$$

$$\vec{w}_{MAP} = \frac{1}{\sigma^2} \left(\alpha \mathbb{I} + \frac{1}{\sigma^2} X^T X \right)^{-1} X^T \vec{y} = \left(\frac{\alpha}{\sigma^2} \mathbb{I} + X^T X \right)^{-1} X^T \vec{y}$$

Therefore $\vec{w}_{MAP} = \left(\lambda \mathbb{I} + X^T X \right)^{-1} X^T \vec{y}$ where $\lambda = \frac{\alpha}{\sigma^2}$

$\left(\underset{\vec{w}}{\operatorname{argmax}} p(\vec{w} | \vec{m}_N, S_N) \right) \Rightarrow \vec{w} = \vec{m}_N$

We can also reason the following way:

$$p(\vec{w} | \vec{y}) = p(\vec{y} | \vec{w}) p(\vec{w}) \Rightarrow \ln p(\vec{w} | \vec{y}) = \ln p(\vec{y} | \vec{w}) + \ln p(\vec{w})$$

$$\ln p(\vec{y} | \vec{w}) = -\frac{1}{2\sigma^2} \sum_i (y_i - \vec{w}^T \vec{x}_i)^2 + \text{const}$$

$$\ln p(\vec{w}) = -\frac{1}{2} \alpha \vec{w}^T \vec{w}$$

$$\vec{w}_{MAP} = \underset{\vec{w}}{\operatorname{argmax}} [\ln p(\vec{w} | \vec{y})]$$

$$\frac{\partial \ln p(\vec{w} | \vec{y})}{\partial \vec{w}} = \frac{1}{\sigma^2} \sum_i (y_i - \vec{w}^T \vec{x}_i) \vec{x}_i - \alpha \vec{w} = 0$$

$$\vec{w}^T (\sigma^2 \alpha \mathbb{I} + X^T X) = X^T \vec{y} \Rightarrow \vec{w}^T = \left(\lambda \mathbb{I} + X^T X \right)^{-1} X^T \vec{y}$$

3. Predictive distribution

It is not only important to find \vec{w}_{ML} or $p(\vec{w})$, but we also are interested in making predictions of y_i given a particular \vec{x}_i .

We can do: $p(y_i | \vec{x}_i) = N(y_i | \vec{w}_{ML}^T \vec{x}_i, \sigma^2)$

or we can take a Bayesian approach that will ~~instead~~ furthermore prevent overfitting:

~~$p(y_i | \vec{y}, \vec{x}_i, X)$~~

$$p(y_i | \vec{x}_i, \vec{y}, X, \alpha, \sigma^2) = \int p(y | \vec{x}_i, \vec{w}, \sigma^2) p(\vec{w} | \vec{y}, X, \alpha, \sigma^2) d\vec{w}$$

where $p(y | \vec{x}_i, \vec{w}, \sigma^2) = N(y | \vec{w}^T \vec{x}_i, \sigma^2)$

$$p(\vec{w} | \vec{y}, X, \alpha, \sigma^2) = N(\vec{w} | \vec{m}_N, S_N)$$

$$\vec{m}_N \approx S_N^{-1} \vec{S}$$

$$\begin{cases} \vec{m}_N = \frac{1}{\sigma^2} S_N X^T \vec{y} \\ S_N^{-1} = \alpha I + \frac{1}{\sigma^2} X^T X \end{cases}$$

Exercise

$$\Rightarrow p(y | \vec{x}_i, \vec{y}, X, \alpha, \sigma^2) = N(y | \vec{m}_N^T \vec{x}_i, \sigma_{N(\vec{x})}^2)$$

where $\sigma_{N(\vec{x})}^2 = \sigma^2 + \vec{x}_i^T S_N \vec{x}_i$

$$\sigma_N^2 = \sigma^2 + \sigma^2 \mathbf{x}^T \left(\mathbf{1} + \frac{\mathbf{x} \mathbf{x}^T}{\sigma^2} \right)^{-1} \mathbf{x}$$

if $N \rightarrow \infty$

$$\begin{aligned} \sigma_N^2 &= \sigma^2 + \sigma^2 \mathbf{x}^T \left(\mathbf{x}^T \mathbf{x} \right)^{-1} \mathbf{x} \\ &= \sigma^2 + \sigma^2 \mathbf{x}^T \mathbf{x}^{-1} \mathbf{x}^{T-1} \mathbf{x} \\ &= \sigma^2 \left(\mathbf{1} + \mathbf{x}^T \mathbf{x}^{-1} \mathbf{x}^{T-1} \mathbf{x} \right) \end{aligned}$$

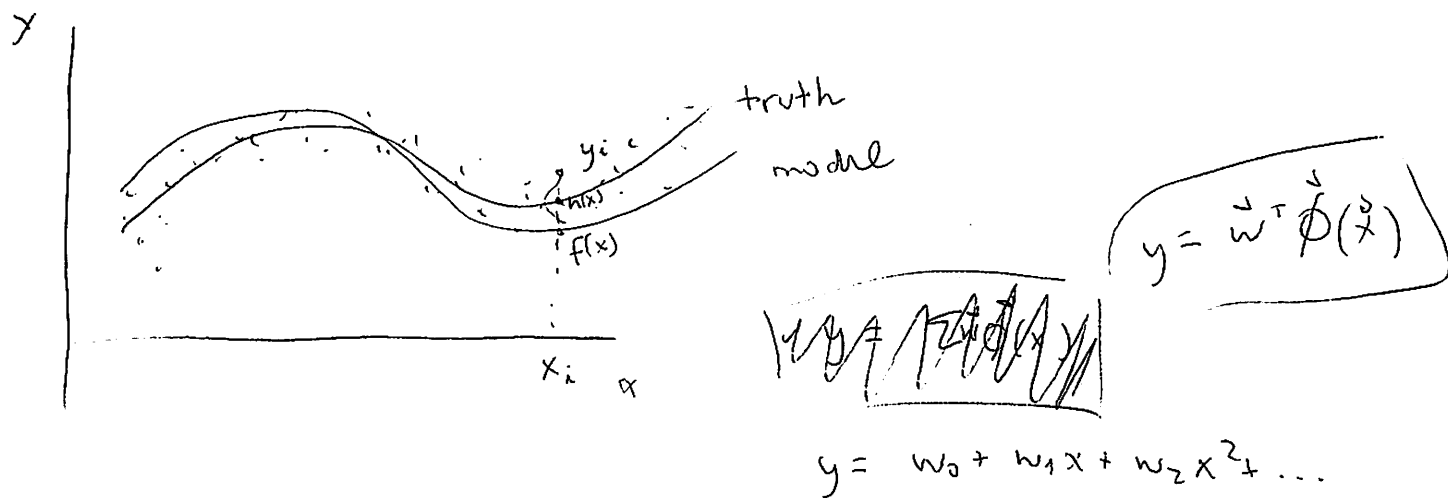
if $N \rightarrow \infty$

$$\begin{aligned} \sigma_N(\mathbf{x}) &= \sigma^2 + \sigma^2 \mathbf{x}^T \left(\mathbf{x}^T \mathbf{x} \right)^{-1} \mathbf{x} \\ &= \sigma^2 \left(1 + \mathbf{x}^T \left(\mathbf{x}^T \mathbf{x} \right)^{-1} \mathbf{x} \right) = \sigma^2 \\ &\quad \downarrow \\ &\quad \text{if } N \rightarrow \infty \\ &\quad \left(\mathbf{x}^T \mathbf{x} \right)^{-1} \rightarrow 0 \end{aligned}$$

The more data the more precise is our estimate.
We can never break the ceiling imposed by the noise of the data itself.

4. Bias-Variance Trade-off

- Let's assume we want to use a more flexible model, not just linear regression:



- If we let \vec{w} full flexibility we can end-up overfitting.
What is the best strength for the regularization?

Let's define:

$$\begin{cases} h(x) = \int y P_t(y|x) dy = \mathbb{E}_t[y|x] \leftarrow \text{True} \\ f(x) = \int y P_m(y|x) dy = \mathbb{E}_m[y|x] \leftarrow \text{model} \end{cases}$$

$$\begin{aligned} \mathbb{E}[\mathcal{L}] &= \mathbb{E}_{x,y}[(y - f(x))^2] = \mathbb{E}_{x,y}[(y - h(x) + h(x) - f(x))^2] \\ &= \mathbb{E}_{x,y}[(y - h(x))^2] + \underbrace{\mathbb{E}_x[(h(x) - f(x))^2]} \end{aligned}$$

Exercise

• Let's focus now on the term $\mathbb{E}_x[(f(x) - h(x))^2]$

Important to note $f(x) = f(x; D)$ where D is the dataset we used to learn our model

$$f(x) = h(x)$$

$$(f(x; D) - h(x))^2 = (f(x; D) - \mathbb{E}_D[f(x; D)] + \mathbb{E}_D[f(x; D)] - h(x))^2$$

$$\Rightarrow \mathbb{E}_D[(f(x; D) - h(x))^2] = \{\mathbb{E}_D[f(x; D)] - h(x)\}^2 + \mathbb{E}_D[(f(x; D) - \mathbb{E}_D[f(x; D)])^2]$$

Exercise

Therefore: $\mathbb{E}_{x,y}[\mathcal{L}] = \mathbb{E}_{x,y}[(y - f(x))^2] =$

$$= \{\mathbb{E}_D[f(x; D)] - h(x)\}^2 \quad \text{Bias}^2$$

$$+ \mathbb{E}_D[(f(x; D) - \mathbb{E}_D[f(x; D)])^2] \quad \text{Variance}$$

$$+ \mathbb{E}_{x,y}[(h(x) - y)^2] \quad \text{Noise}$$

