

Lecture 16

Linear regression

Consider $p(y|\vec{x}, \vec{\theta}) = N(y | \underbrace{\vec{w} \cdot \vec{x}}_{D\text{-dim vector}} + w_0, \sigma^2)$,
 where $\vec{\theta} = (\vec{w}, w_0, \sigma^2)$.

Write $\vec{x} = (1, x_1, \dots, x_D)$, absorb w_0 into \vec{w} : $\vec{w} \rightarrow (w_0, \vec{w})$.

Features: $p(y|\vec{x}, \vec{\theta}) = N(y | \underbrace{\vec{w} \cdot \vec{f}(\vec{x})}_{M\text{-dim vector}}, \sigma^2)$

Multi-D version:

$$p(\vec{y}|\vec{x}, \vec{w}) = \prod_{j=1}^M N(y_j | \underbrace{\vec{w}_j \cdot \vec{x}}_{\substack{\uparrow \\ D+1 \\ \text{dim}}}, \sigma_j^2)$$

factorizes into M 1D problems

with correlations,

$$p(\vec{y}|\vec{x}, \vec{w}) = N(\vec{y} | \underbrace{\vec{w} \cdot \vec{x}}_{M \times (D+1)}, \Sigma)$$

\uparrow
covariance matrix

Next, consider the 1D case:

$$-\mathcal{J}(\vec{w}, \sigma^2) = \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \vec{w} \cdot \vec{x}_n)^2 +$$

$$+ \frac{N}{2} \log(2\pi\sigma^2).$$

\equiv

Focus first on [i.e., infer \vec{w} first,
 σ^2 second]

$$RSS(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (y_n - \vec{w} \cdot \vec{x}_n)^2 \quad \text{①}$$

↑
 residual sum
 of squares ② $\frac{1}{2} (\vec{y} - X\vec{w}) \cdot (\vec{y} - X\vec{w})$.
 \uparrow
 dim = N

$X \vec{w}_{N,D+1}^{D+1}$ = N -dim vector.

data matrix, each row is
 an input pattern (augmented with 1)

Note that

$$\underbrace{\frac{\partial}{\partial w_i} RSS(\vec{w})}_{D+1 \text{-dim vector}} = - \sum_{n=1}^N \underbrace{(y_n - \vec{w} \cdot \vec{x}_n)}_{\vec{y} - X\vec{w}} \underbrace{x_{n,i}}_{N\text{-dim}}.$$

$\downarrow \quad \downarrow \quad \downarrow$
 $X^T_{D+1,N}$

Then

$$+ \underbrace{\frac{\partial}{\partial \vec{w}} RSS(\vec{w})}_{\Leftrightarrow g(\vec{w})} = \underbrace{X^T X \vec{w}}_{(D+1) \times (D+1)} - \underbrace{X^T \vec{y}}_{} = 0, \text{ s.t.}$$

$$\vec{w} = \underbrace{(X^T X)^{-1}}_{\text{Moore-Penrose}} \underbrace{X^T \vec{y}}_{\text{pseudoinverse}}. \quad (*)$$

Moreover, $\frac{\partial^2}{\partial w_i \partial w_j} RSS(\vec{w}) = \sum_n x_{n,i} x_{n,j}$,

$$H(\vec{w}) = \underbrace{X^T X}_{\rightarrow X^T X}$$

Now, for $\forall \vec{v} > 0$ we have

$$\vec{v}^\top (X^\top X) \vec{v} = (X\vec{v})^\top (X\vec{v}) > 0 \text{ if } X \text{ is full-rank.}$$

So, $H(\vec{w})$ is pos.-def., and the $\vec{g}(\vec{w})=0$ solution is unique and is given by Eq. (*).

Typically, $N \gg D \Rightarrow$ the X matrix is 'tall and skinny'. This is an overdetermined system of linear equations.

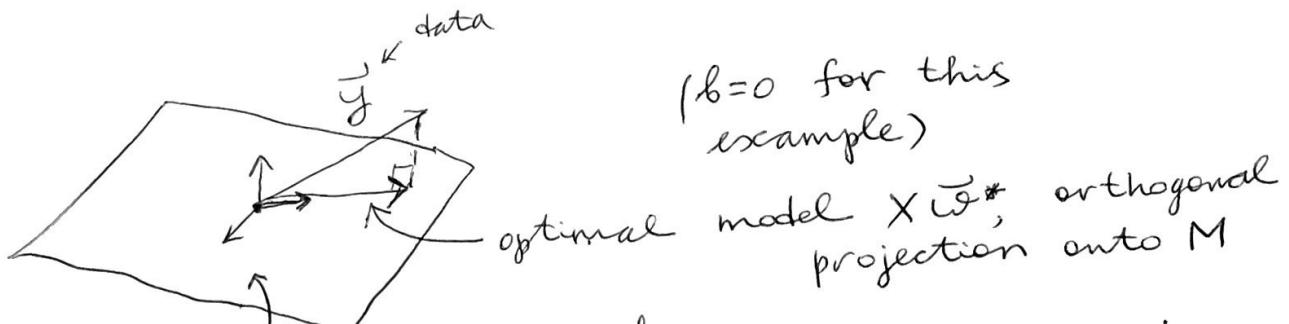
Geometrically, we seek to represent an N -dim vector $\vec{y} = \underline{y_1, y_2, \dots, y_N}$ as

a linear combination of $D+1$ N -dim vectors $\vec{x}_1^c, \dots, \vec{x}_{D+1}^c \in$ columns of X .

These vectors span $(D+1)$ -dim subspace of the N -dim space if X is full rank.

In our model, we consider $X\vec{w} = w_1 \vec{x}_1^c + \dots + w_{D+1} \vec{x}_{D+1}^c$ as our prediction.

We want to minimize the magnitude of the difference between \vec{y} & $X\vec{w}$.



$M = \text{subspace spanned by } \bar{x}_1^c, \dots, \bar{x}_{D+1}^c, \text{ embedded into } \mathbb{R}^N$

Clearly, $\bar{x}_i^c \cdot (\bar{y} - X\bar{w}^*) = 0, \forall i = 1, \dots, D+1$

Then $\underbrace{X^T}_{D+1, N} (\bar{y} - X\bar{w}^*) = \vec{0}_{D+1}, \text{ or}$

$$(X^T X) \bar{w}^* = X^T \bar{y},$$

$$\bar{w}^* = (X^T X)^{-1} X^T \bar{y}, \text{ same as Eq. (*) above}$$

If we denote $\hat{y} = X\bar{w}^*$,

$$\hat{y} = \underbrace{X(X^T X)^{-1} X^T \bar{y}}_{\text{projection operator}}$$

Numerically, using Eq. (*) is to invert $X^T X$ or compute $(X^T X)^{-1} X$ using SVD or other methods.

Weighted least squares

Consider $p(y|\vec{x}, \vec{\omega}) = N(y|\vec{\omega} \cdot \vec{x}, \underbrace{\sigma^2(\vec{x})}_{\text{variance}})$
 variance depends on input \vec{x}

Then $p(\vec{y}|\vec{x}, \vec{\omega}) = N(\vec{y}|\vec{X}\vec{\omega}, \Lambda^{-1})$, where

$$\vec{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \quad \Lambda = \text{diag} \left[\frac{1}{\sigma^2(\vec{x}_1)}, \dots, \frac{1}{\sigma^2(\vec{x}_N)} \right].$$

$$RSS(\vec{\omega}) = \underbrace{\frac{1}{2} \sum_{n=1}^N (y_n - \vec{\omega} \cdot \vec{x}_n)^2}_{\text{weighted sum of squares}} \frac{1}{\sigma^2(\vec{x}_n)}.$$

If the weights are known, we obtain:

$$\frac{\partial}{\partial \omega_i} RSS(\vec{\omega}) = - \sum_n (y_n - \vec{\omega} \cdot \vec{x}_n) \frac{x_{n,i}}{\sigma^2(\vec{x}_n)}.$$

$$\text{Then, if } \frac{\partial}{\partial \vec{\omega}} RSS(\vec{\omega}) = 0 = (\vec{X}^\top \Lambda \vec{X}) \vec{\omega} - \vec{X}^\top \Lambda \vec{y}, \text{ or}$$

$$\vec{\omega}^* = (\vec{X}^\top \Lambda \vec{X})^{-1} \vec{X}^\top \Lambda \underline{\vec{y}}.$$

Sometimes it makes sense to compute the bias and the D-dim weights separately.

Specifically, consider $\underset{D\text{-dim}}{\stackrel{\omega_0}{\text{RSS}}}(\vec{\omega}) = \frac{1}{2} \sum_{n=1}^N (y_n - b - \vec{\omega} \cdot \vec{x}_n)^2$.

Then $\frac{\partial}{\partial b} \text{RSS}(\vec{\omega}) = -\sum_{n=1}^N (y_n - b - \vec{\omega} \cdot \vec{x}_n) = 0$, or

$$Nb = \underbrace{\sum_n y_n}_{N\langle y \rangle} - \vec{\omega} \cdot \underbrace{\sum_n \vec{x}_n}_{N\langle \vec{x} \rangle},$$

(**) $b = \underbrace{\langle y \rangle}_{\text{ave of data}} - \underbrace{\vec{\omega} \cdot \langle \vec{x} \rangle}_{\substack{\text{ave of each component} \\ \text{of } \vec{x} \text{ over } N \text{ datapoints}}}$

Next, $\text{RSS}(\vec{\omega}) = \frac{1}{2} \sum_{n=1}^N \underbrace{(y_n - \langle y \rangle)}_{y_n^c} - \vec{\omega} \cdot \underbrace{(\vec{x}_n - \langle \vec{x} \rangle)}_{\vec{x}_n^c}$.

Same as before, but with
 D -dim \vec{x} and both \vec{x} & y 'centered'
 by subtracting the mean.

Then $\vec{\omega}_D^* = (\vec{x}_c^T \vec{x}_c)^{-1} \vec{x}_c^T \vec{y}_c$ and

b is found using Eq. (**).

Note that here

$$\vec{\omega}^* = \left[\underbrace{\sum_{n=1}^N (\vec{x}_n - \langle \vec{x} \rangle)(\vec{x}_n - \langle \vec{x} \rangle)^T}_{D \times D \text{ matrix}} \right]^{-1} \times$$

$$\times \underbrace{\sum_{n=1}^N (y_n - \langle y \rangle)(\vec{x}_n - \langle \vec{x} \rangle)^T}_{D\text{-dim vector}}.$$

If $D=1$,

$$\left\{ \begin{array}{l} w_1^* = \frac{\sum_n (x_n - \langle x \rangle)(y_n - \langle y \rangle)}{\sum_n (x_n - \langle x \rangle)(x_n - \langle x_n \rangle)} = \frac{\text{Cov}(x, y)}{\text{Var}(x)}, \\ b = \langle y \rangle - w_1^* \langle x \rangle. \end{array} \right.$$

Recall that

$$-\mathcal{J}(\vec{w}, \sigma^2) = \frac{\text{RSS}(\vec{w})}{\sigma^2} + \frac{N}{2} \log(2\pi\sigma^2).$$

Then

$$\frac{\partial}{\partial \sigma^2} (-\mathcal{J}(\vec{w}, \sigma^2)) = -\frac{\text{RSS}(\vec{w})}{(\sigma^2)^2} + \frac{N}{2\sigma^2} = 0,$$

or

$$\sigma^2 = \frac{2}{N} \text{RSS}(\vec{w}) = \frac{1}{N} \sum_{n=1}^N (y_n - \vec{w} \cdot \vec{x}_n)^2.$$


 MSE of the residuals