# Assumptions
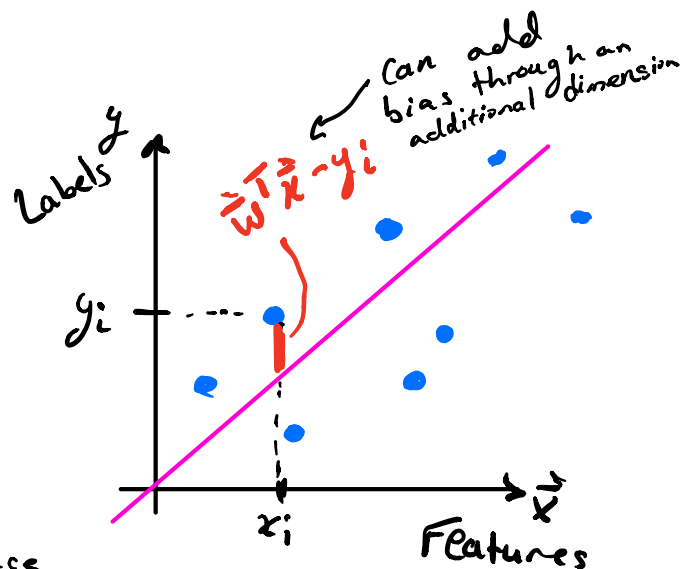
Data: $y_i \in \mathbb{R}$

Model: $y_i = \bar{w}^T \bar{x}_i + \varepsilon_i$, where $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\Rightarrow y | \bar{x}_i \sim \mathcal{N}(\bar{w}^T \bar{x}_i, \sigma^2)$$

$$\Rightarrow \mathbb{P}(y | \bar{x}_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\bar{x}_i^T \bar{w} - y_i)^2}{2\sigma^2}}$$

In words, this means we can assume the data is drawn from a "line" $\vec{w}^T \vec{x}$ through the origin. For each data point w/ features $\vec{x}_i$, the label $y$ is drawn from a Gaussian w/ mean $\vec{w}^T \vec{x}_i$ and variance $\sigma^2$.

Our task is to estimate the slope $\vec{w}$ from the data.

# Estimating with MLE

$$\vec{w} = \underset{w}{\arg\max}\ \mathbb{P}\big((y_1, \vec{x_1}), (y_2, \vec{x_2}), \ldots, (y_n, \vec{x_n}) \mid \vec{w}\big)$$

$$= \underset{w}{\arg\max}\ \prod_{i=1}^{n} \mathbb{P}\big((y_i, \vec{x_i}) \mid w\big) \qquad \textcolor{red}{\text{By independence}}$$

$$= \underset{w}{\arg\max}\ \prod_{i=1}^{n} \mathbb{P}\big(y_i \mid \vec{x_i}, \vec{w}\big)\, \mathbb{P}\big(\vec{x_i} \mid \vec{w}\big)$$

red side note:
$$\textcolor{red}{\mathbb{P}(A \mid B, C)\, \mathbb{P}(B \mid C)}$$
$$\textcolor{red}{= \frac{\mathbb{P}(A \cap B \cap C)}{\mathbb{P}(B \cap C)}\, \frac{\mathbb{P}(B \cap C)}{\mathbb{P}(C)}}$$
$$\textcolor{red}{= \mathbb{P}(A, B \mid C)}$$

$$= \underset{w}{\arg\max}\ \prod_{i=1}^{n} \mathbb{P}\big(y_i \mid \vec{x_i}, \vec{w}\big)\, \mathbb{P}(\vec{x_i}) \qquad \textcolor{red}{\vec{x_i} \text{ independent of } \vec{w}\ \forall i}$$

$$= \underset{w}{\arg\max}\ \prod_{i=1}^{n} \mathbb{P}\big(y_i \mid \vec{x_i}, \vec{w}\big) \qquad \textcolor{red}{\mathbb{P}(\vec{x_i}) \text{ is a constant}}$$

$$= \underset{w}{\arg\max}\ \sum_{i=1}^{n} \log\big(\mathbb{P}(y_i \mid \vec{x_i}, \vec{w})\big) \qquad \textcolor{red}{\log \text{ is a monotonic function}}$$

$$= \underset{w}{\arg\max}\ \sum_{i=1}^{n} \log\left( \frac{1}{\sqrt{2\pi\sigma^2}}\, e^{-\frac{(\vec{x_i}^T \vec{w} - y_i)^2}{2\sigma^2}} \right)$$

$$= \underset{w}{\arg\max}\ \sum_{i=1}^{n} \log\left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \log\left( e^{-\frac{(\vec{x_i}^T \vec{w} - y_i)^2}{2\sigma^2}} \right)$$

$$= \underset{w}{\arg\max}\ \sum_{i=1}^{n} -\frac{(\vec{x_i}^T \vec{w} - y_i)^2}{2\sigma^2}$$

$$\textcolor{magenta}{\textbf{ALWAYS} \text{ wanna minimize}}$$

$$= \underset{w}{\arg\min}\ \frac{1}{n} \sum_{i=1}^{n} (\vec{x_i}^T \vec{w} - y_i)^2 \quad \textcolor{purple}{\Big\}\ \text{Ordinary Least Squares}}$$

$$\textcolor{purple}{\text{optimize w/ gradient descent}}$$

$$\textcolor{red}{\frac{1}{2\sigma^2} \text{ is a constant which doesn't change } \vec{w}. \text{ dividing by } n \text{ instead makes the loss more tractable.}}$$

**Closed Form:** $\vec{w} = (XX^T)^{-1} X \vec{y}^T$ ; where $X = [x_1, \ldots, x_n]$, $y = [y_1, \ldots, y_n]$

# Estimating with MAP

Additional Model Assumptions:

$$P(\vec{w}) = \frac{1}{\sqrt{2\pi\tau^2}} e^{-\frac{\vec{w}^T\vec{w}}{2\tau^2}}$$

MAP maximizes $w$ from $P(w \mid \text{dataset})$

Note: $D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \ldots, (\vec{x}_n, y_n)\}$

$$\vec{w} = \underset{w}{\arg\max} \; P(\vec{w} \mid D)$$

$$= \underset{w}{\arg\max} \; P(\vec{w} \mid (\vec{x}_1, y_1), (\vec{x}_2, y_2), \ldots, (\vec{x}_n, y_n))$$

$$= \underset{w}{\arg\max} \; \frac{P((\vec{x}_1, y_1), (\vec{x}_2, y_2), \ldots, (\vec{x}_n, y_n) \mid \vec{w}) \, P(\vec{w})}{P((\vec{x}_1, y_1), (\vec{x}_2, y_2), \ldots, (\vec{x}_n, y_n))} \quad \leftarrow \text{constant}$$

$$= \underset{w}{\arg\max} \; P((\vec{x}_1, y_1), (\vec{x}_2, y_2), \ldots, (\vec{x}_n, y_n) \mid \vec{w}) \, P(\vec{w})$$

$$= \underset{w}{\arg\max} \; \prod_{i=1}^{n} \left[ P((x_i, y_i) \mid \vec{w}) \right] P(\vec{w})$$

$$= \underset{w}{\arg\max} \; \prod_{i=1}^{n} \left[ P(y_i \mid \vec{x}_i, \vec{w}) \, P(\vec{x}_i \mid \vec{w}) \right] P(\vec{w})$$

$$= \underset{w}{\arg\max} \; \prod_{i=1}^{n} \left[ P(y_i \mid \vec{x}_i, \vec{w}) \, P(\vec{x}_i) \right] P(\vec{w})$$

$$= \underset{w}{\arg\max} \; \prod_{i=1}^{n} \left[ P(y_i \mid \vec{x}_i, \vec{w}) \right] P(\vec{w})$$

$$= \underset{w}{\arg\max} \; \sum_{i=1}^{n} \log\left( P(y_i \mid \vec{x}_i, \vec{w}) \right) + \log\left( P(\vec{w}) \right)$$

$$= \underset{w}{\arg\max} \; \sum_{i=1}^{n} \log\left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\vec{x}_i^T\vec{w} - y_i)^2}{2\sigma^2}} \right) + \log\left( \frac{1}{\sqrt{2\pi\tau^2}} e^{-\frac{\vec{w}^T\vec{w}}{2\tau^2}} \right)$$

$$= \underset{w}{\text{argmax}} \sum_{i=1}^{n} \left[ -\frac{1}{2\sigma^2} (\vec{x}_i^T \vec{w} - y_i)^2 - \frac{\vec{w}^T \vec{w}}{2\tau^2} \right]$$

$$= \underset{w}{\text{argmin}} \sum_{i=1}^{n} (\vec{x}_i^T \vec{w} - y_i)^2 - \frac{n\sigma^2}{\tau^2} \vec{w}^T \vec{w}$$

$$= \underset{w}{\text{argmin}} \frac{1}{n} \sum_{i=1}^{n} (\vec{x}_i^T \vec{w} - y_i)^2 - \lambda \vec{w}^T \vec{w} \quad ; \quad \lambda = \frac{\sigma^2}{n\tau^2}$$

<span style="color:red">↑</span>

<span style="color:red">Known as ridge regression. Differs from OLS in that there's $L_2$ regularization</span>

Closed form: $\vec{w} = (XX^T - \lambda I)^{-1} X \vec{y}^T$

where $X = [x_1, \dots, x_n], y = [y_1, \dots, y_n]$