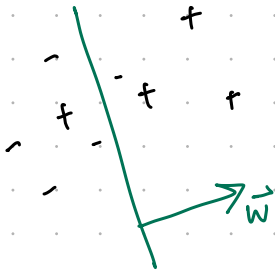


$$P(y | \vec{x}) = \frac{1}{1 + e^{-y \vec{w}^T \vec{x}}} \quad \text{for } y \in \{-1, +1\}$$



- By viewing \vec{w} as parameter

$$\vec{w}_{MLE} = \underset{\vec{w}}{\operatorname{argmin}} \underbrace{\sum_{i=1}^n \log(1 + e^{-y_i \vec{w}^T \vec{x}_i})}_{\ell(\vec{w}) = \sum_{i=1}^n L(x_i, y_i, \vec{w})}$$

- Claim that $\ell(\vec{w})$ is continuous, differentiable, and convex.
(logistic loss function) ✓ ✓ ✓

- The term $\vec{w}^T \vec{x}_i$ is the prediction of \vec{x}_i

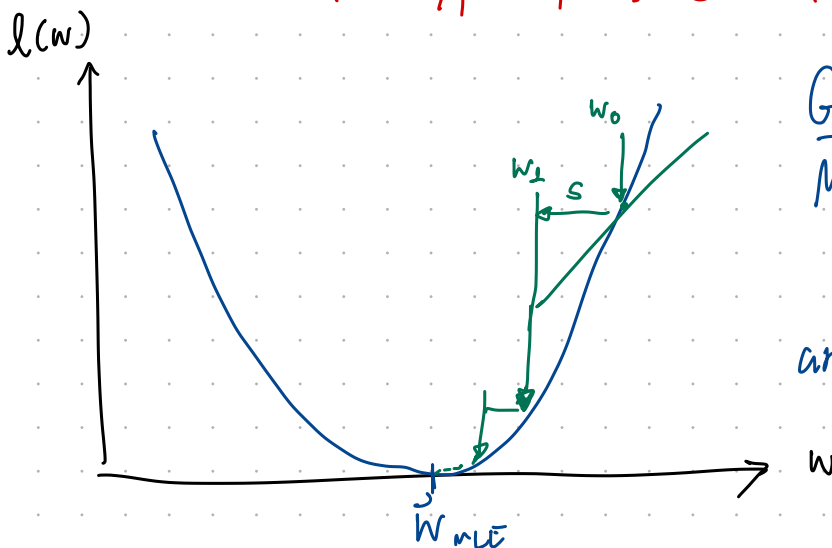
$$h(\vec{x}_i) = \operatorname{sign}(\vec{w}^T \vec{x}_i)$$

- If the prediction is correct, then $h(\vec{x}_i) = y_i$

$$\Rightarrow -y_i \vec{w}^T \vec{x}_i < 0 \Rightarrow L(x_i, y_i, \vec{w}) \rightarrow 0$$

- If the prediction is wrong, then $h(\vec{x}_i) \neq y_i$

$$\Rightarrow -y_i \vec{w}^T \vec{x}_i > 0 \Rightarrow L(x_i, y_i, \vec{w}) \rightarrow y_i \vec{w}^T \vec{x}_i$$



Gradient Descent

Make a progress with step

$$\vec{s} = -\alpha \nabla \ell(\vec{w})$$

and repeat the process until
converged

We need to ensure that $\ell(w_{j+1}) < \ell(w_j)$ at each step j

Taylor Expansion

- If $\|\vec{s}\|_2$ is small (meaning $\vec{w} + \vec{s}$ is close to \vec{w}), then the following holds

$$l(\vec{w} + \vec{s}) \approx l(\vec{w}) + \nabla l(\vec{w})^T \vec{s}$$

- With $\vec{s} = -\alpha \nabla l(\vec{w})$, where $\alpha > 0$, we have that

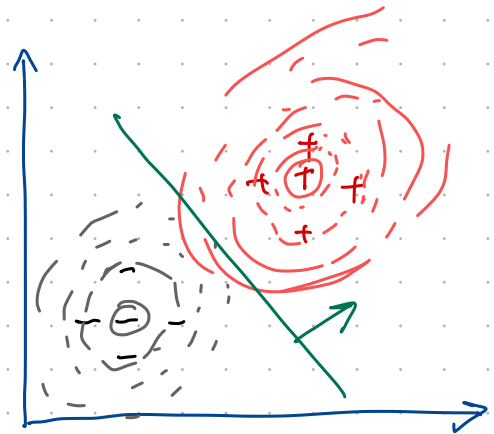
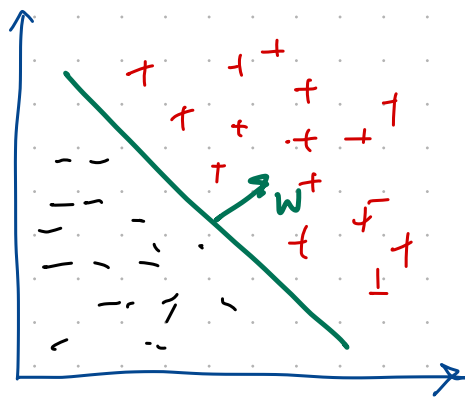
$$\begin{aligned} l(\vec{w} + \vec{s}) &\approx l(\vec{w}) + \nabla l(\vec{w})^T (-\alpha \nabla l(\vec{w})) \\ &= l(\vec{w}) + (-\alpha (\nabla l(\vec{w})^T) \nabla l(\vec{w})) \\ &= l(\vec{w}) + (-\alpha \underbrace{\|\nabla l(\vec{w})\|_2}_{>0})^2 \\ &< l(\vec{w}) \end{aligned}$$

What's to pick up

- Logistic loss function $l(\vec{w}) = \sum_{i=1}^n \log(1 + e^{-y_i \vec{w}^T \vec{x}_i})$ also measures how well the linear classifier performs on our data. Ideally, we want $l(\vec{w}) \rightarrow 0$.
- The Gradient Descent is a hill climbing algorithm that finds an optimal point of several loss function
- Logistic regression is an ML algorithm for binary classification that gives a linear classifier (discriminative counterpart of naïve Bayes)
- No assumption about model's distributions made for Logistic regression

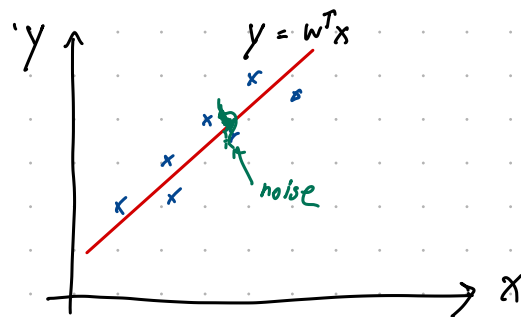
Naive Bayes vs Logistic Regression

- With little data and if the modeling distribution is appropriate, Naive Bayes tends to outperform Logistic Regression
- As data sets become larger, Logistic Regression often outperforms Naive Bayes, which suffers from the fact that the modeling assumptions made for $P(X|y)$ are probably not exactly correct



Linear Regression

- Assume $y_i \in \mathbb{R}$
- Assume $y_i = \underbrace{\vec{w}^T \vec{x}_i}_{\text{line}} + \underbrace{\varepsilon_i}_{\text{noise}}$,
where $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$



$$\Leftrightarrow y_i \sim \mathcal{N}(\vec{w}^T \vec{x}_i, \sigma^2) \Rightarrow P(y_i | \vec{x}_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\vec{w}^T \vec{x}_i - y_i)^2}{2\sigma^2}}$$

The goal is to estimate the vector \vec{w} that defines the slope.

Estimate \vec{w} with MLE

$$\begin{aligned}\vec{w}_{MLE} &= \arg \max_{\vec{w}} \prod_{i=1}^d P(y_i | \vec{x}_i; \vec{w}) \\&= \arg \max_{\vec{w}} \sum_{i=1}^d \log(P(y_i | \vec{x}_i; \vec{w})) \\&= \arg \max_{\vec{w}} \sum_{i=1}^d \left[\log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \log\left(e^{-\frac{(\vec{w}^T \vec{x}_i - y_i)^2}{2\sigma^2}}\right) \right] \\&= \arg \max_{\vec{w}} -\frac{1}{2\sigma^2} \sum_{i=1}^n (\vec{w}^T \vec{x}_i - y_i)^2 \\&= \arg \min_{\vec{w}} \underbrace{\frac{1}{n} \sum_{i=1}^n (\vec{w}^T \vec{x}_i - y_i)^2}_{l(\vec{w})}\end{aligned}$$

** Remark: Here, we are minimizing the square loss function $l(\vec{w})$

Ways to find \vec{w}_{MLE} : ① "Gradient Descent" since the square loss is continuous, differentiable, and convex (parabola)

② Compute via a closed-form solution of the square loss function

Closed-form solution of \vec{w}_{MLE}

$$\vec{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}, \text{ where } \mathbf{X} = \begin{bmatrix} \leftarrow \vec{x}_1^T \rightarrow \\ \leftarrow \vec{x}_2^T \rightarrow \\ \vdots \\ \leftarrow \vec{x}_n^T \rightarrow \end{bmatrix} \text{ and } \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\begin{aligned}
 X \vec{w} - \vec{y} &= \begin{bmatrix} \vec{x}_1^T \\ \vec{x}_2^T \\ \vdots \\ \vec{x}_n^T \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \\
 &= \begin{bmatrix} \vec{w}^T \vec{x}_1 \\ \vec{w}^T \vec{x}_2 \\ \vdots \\ \vec{w}^T \vec{x}_n \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \vec{w}^T \vec{x}_1 - y_1 \\ \vec{w}^T \vec{x}_2 - y_2 \\ \vdots \\ \vec{w}^T \vec{x}_n - y_n \end{bmatrix}
 \end{aligned}$$

$$l(\vec{w}) = (X\vec{w} - \vec{y})^T (X\vec{w} - \vec{y}) = \sum_{i=1}^n (\vec{w}^T \vec{x}_i - y_i)^2$$

$$\begin{aligned}
 \nabla_{\vec{w}} l(\vec{w}) &= \nabla_{\vec{w}} (X\vec{w} - \vec{y})^T (X\vec{w} - \vec{y}) \\
 &= \nabla_{\vec{w}} (w^T X^T X \vec{w} - y^T X \vec{w} - X^T w^T \vec{y} + \cancel{\vec{y}^T \vec{y}}) \\
 &= 2 X^T X \vec{w} - y^T X - X^T y
 \end{aligned}$$

$$\boxed{\nabla_{\vec{w}} l(\vec{w}) = 2 X^T X \vec{w} - 2 X^T y}$$

To find \vec{w}_{MLE} , $\nabla_{\vec{w}} l(\vec{w}) \stackrel{\text{set}}{=} 0$

$$\Rightarrow \cancel{2} X^T X \vec{w}_{MLE} = \cancel{2} X^T y$$

$$\Rightarrow X^T X \vec{w}_{MLE} = X^T y$$

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

← Closed-Form solution
(complexity $O(d^3)$)

*** If d is large, the solution is not possible to compute.