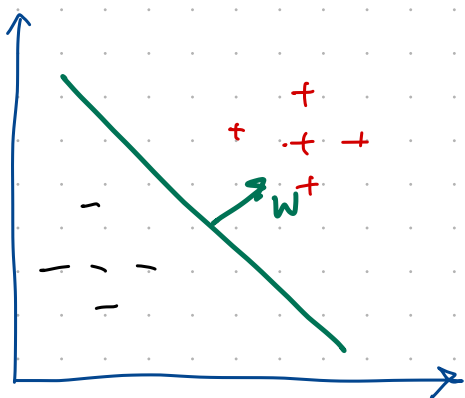


$$P(y|x) \propto P(x|y) P(y)$$

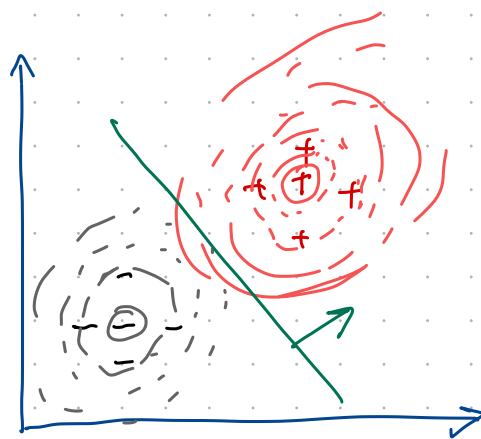
Types of ML algorithms:

- Generative algorithms: try to model $P(x|y)$ and $P(y)$.
- Discriminative algorithms: try to model $P(y|x)$ directly.



The perceptron is a discriminative algorithm since it models

$$P(y = +1 | \vec{x}) = \begin{cases} 1 & \text{if } \vec{w}^T \vec{x} > 0 \\ 0 & \text{o.t.w.} \end{cases}$$



Naive Bayes is a generative algorithm since it models


$P(\vec{x} | y = +1)$ with some explicit modeling distribution (e.g. Gaussian).

- From our last lecture, we have shown that if the data is multinomial features, then Naive Bayes is a linear classifier

$$h(\vec{x}) = +1 \quad \text{iff} \quad \underbrace{\vec{w}^T \vec{x}}_{\vec{w}^T \vec{x}} + b > 0 \quad \text{for specific vector } \vec{w} \quad \text{and scalar } b \text{ and } y \in \{-1, +1\}$$

- With the similar derivation for the case of Gaussian naive Bayes, we can also derive

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



\vec{x}_1
 $P(y=+1|\vec{x}_1) \approx 1$
 \vec{x}_2
 $P(y=+1|\vec{x}_2) \approx 0$

*** We have a beautiful closed-form solution for $P(y|\vec{x})$ where the vector \vec{w} can be found by fitting the parameter regarding Gaussian distribution (our modeling assumption)

Two brilliant ideas: (1) It is nice to estimate \vec{w} directly
 (2) It is nice to set the sigmoid function for modelling $P(y|\vec{x})$

Logistic Regression

- Discriminative counterpart of naive Bayes.
- The assumption is the following

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

- Here, we use the sigmoid function to model $P(y|\vec{x})$ where the \vec{w} is the parameter that we need to estimate from data

Estimating \vec{w}

- MLE: We choose \vec{w}_{MLE} that maximize the conditional likelihood

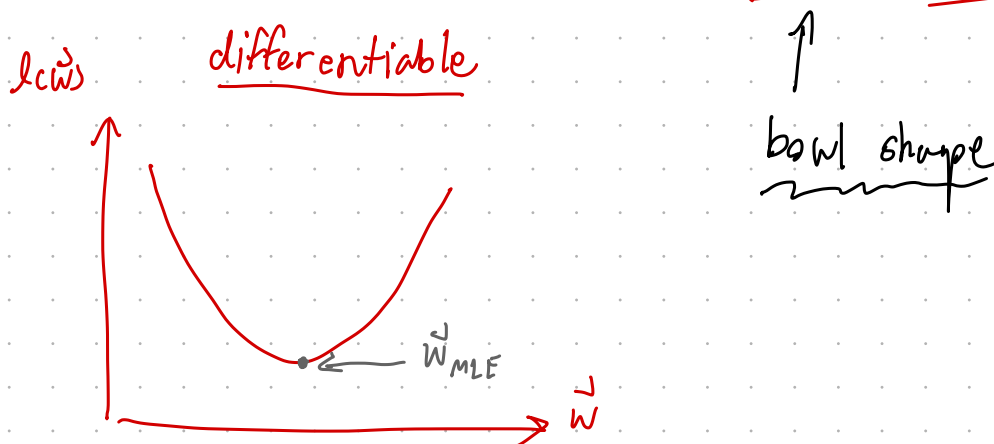
$$\begin{aligned}\vec{w}_{MLE} &= \underset{\vec{w}}{\operatorname{argmax}} \prod_{i=1}^n p(y_i | \vec{x}_i, \vec{w}) \\ &= \underset{\vec{w}}{\operatorname{argmax}} \sum_{i=1}^n \log \left(\frac{1}{1 + e^{-y_i \vec{w}^T \vec{x}_i}} \right) \\ &= \underset{\vec{w}}{\operatorname{argmax}} - \sum_{i=1}^n \log (1 + e^{-y_i \vec{w}^T \vec{x}_i}) \\ &= \underset{\vec{w}}{\operatorname{argmin}} \sum_{i=1}^n \log (1 + e^{-y_i \vec{w}^T \vec{x}_i})\end{aligned}$$

- Note that there is no closed-form solution to \vec{w}_{MLE} .

- We will use Gradient Descent on the negative log likelihood

$$l(\vec{w}) = \sum_{i=1}^n \log (1 + e^{-y_i \vec{w}^T \vec{x}_i}).$$

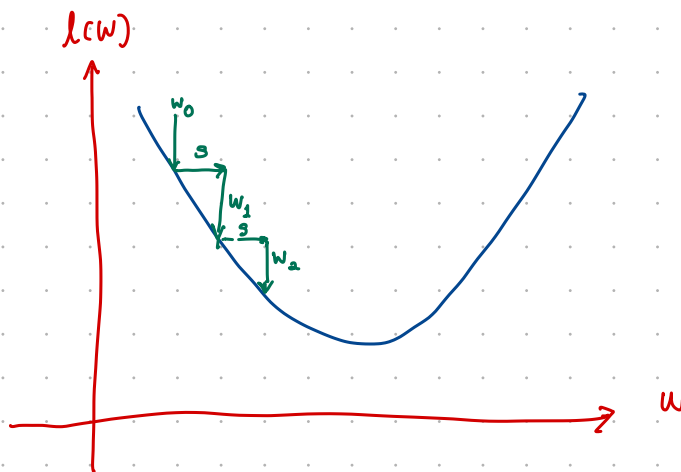
Observation: The function $l(\vec{w})$ is convex, continuous, and differentiable



Goal: Find \vec{w} that minimizes $l(\vec{w})$

Hill-climbing algorithm

1. Initialize \vec{w}_0
2. While $\|\vec{w}_{t+1} - \vec{w}_t\|_2 \geq \epsilon$
3. $\vec{w}_{t+1} = \vec{w}_t + \vec{s}$



The problem is how to define \vec{s} ?

Trick: Taylor's Approximation

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

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

Gradient Descent: First-order approximation

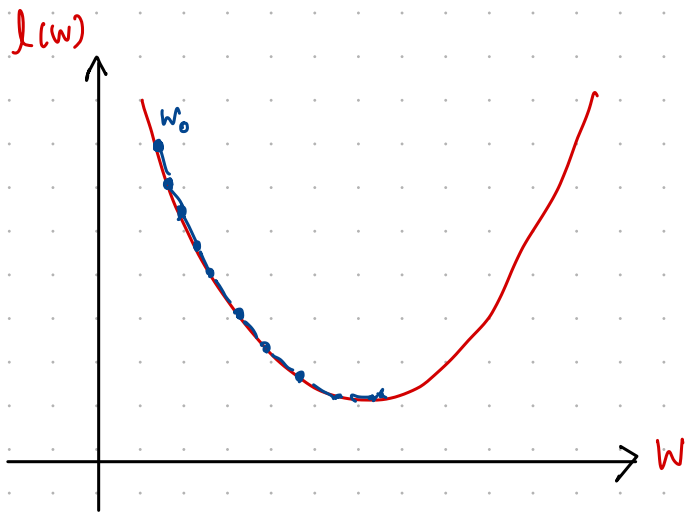
- We assume that the function l around \vec{w} is linear.
- We wish to take step \vec{s} where $l(\vec{w}) > l(\vec{w} + \vec{s})$.
- In gradient descent, we set

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

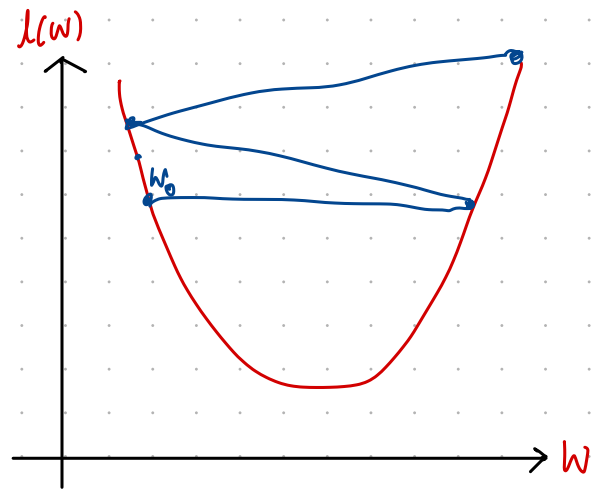
where $\alpha > 0$ is the learning rate.

$$\Rightarrow l(\vec{w} + \vec{s}) = l(\underbrace{\vec{w} + (-\alpha \nabla l(\vec{w}))}_{\text{after one update}}) \approx l(\vec{w}) - \underbrace{\alpha (\nabla l(\vec{w}))^T (\nabla l(\vec{w}))}_{\substack{\|\nabla l(\vec{w})\|_2 > 0 \\ < l(\vec{w})}}$$

Setting the learning rate α



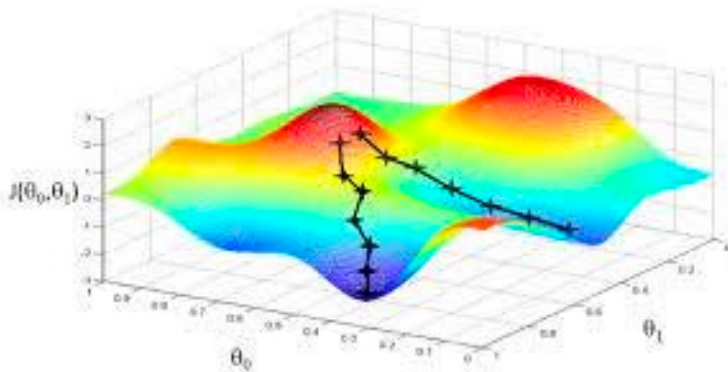
If α is too small,
the algorithm will converge
very slowly.



If α is too large,
the algorithm will never
converge

- A safe choice is to set $\alpha = \frac{t_0}{t_1}$, which guarantees that it will eventually become small enough to converge (For any $t_0 > 0$).

Remark: If the function is not convex, it may converge at local optimum.



$$P(y=1 | \vec{x}) = \frac{P(\vec{x} | y=+1) \times P(y=+1)}{P(\vec{x})}$$

$$= \frac{P(\vec{x} | y=+1) \times P(y=+1)}{P(\vec{x} | y=+1) P(y=+1) + P(\vec{x} | y=-1) P(y=-1)}$$

$$= \frac{1}{1 + e^{\log\left(\frac{P(\vec{x} | y=-1) P(y=-1)}{P(\vec{x} | y=+1) P(y=+1)}\right)}}$$

$$= \frac{1}{1 + e^{\left[\log\left(\frac{P(\vec{x} | y=-1)}{P(\vec{x} | y=+1)}\right) + \log\left(\frac{P(y=-1)}{P(y=+1)}\right)\right]}}$$

$$= \frac{1}{1 + e^{\left[\log\left(\frac{\prod_{i=1}^d P(x_i | y=-1)}{\prod_{i=1}^d P(x_i | y=+1)}\right) + \log\left(\frac{P(y=-1)}{P(y=+1)}\right)\right]}} = \frac{1}{1 + e^{\left(\log\left(\frac{\eta}{1-\eta}\right) + \sum_{i=1}^d \log\left(\frac{P(x_i | y=-1)}{P(x_i | y=+1)}\right)\right)}}$$

$$\begin{aligned} \sum_i \log\left(\frac{P(x_i | y=-1)}{P(x_i | y=+1)}\right) &= \sum_{i=1}^d \log \frac{\frac{1}{\sqrt{2\pi}\sigma_i^2} \exp\left(-\frac{(x_i - \mu_{i,-1})^2}{2\sigma_i^2}\right)}{\frac{1}{\sqrt{2\pi}\sigma_i^2} \exp\left(-\frac{(x_i - \mu_{i,+1})^2}{2\sigma_i^2}\right)} \\ &= \sum_{i=1}^d \log\left(\exp\left(-\frac{(x_i - \mu_{i,-1})^2 + (x_i - \mu_{i,+1})^2}{2\sigma_i^2}\right)\right) \\ &= \sum_{i=1}^d \left(\frac{x_i^2 - 2x_i\mu_{i,+1} + (\mu_{i,+1})^2 - x_i^2 + 2x_i\mu_{i,-1} - (\mu_{i,-1})^2}{2\sigma_i^2}\right) \\ &= \sum_{i=1}^d \left(\frac{2x_i(\mu_{i,-1} - \mu_{i,+1}) + (\mu_{i,+1})^2 - (\mu_{i,-1})^2}{2\sigma_i^2}\right) \\ &= \sum_{i=1}^d \left(\frac{(\mu_{i,-1} - \mu_{i,+1})}{\sigma_i^2} x_i + \frac{(\mu_{i,+1})^2 - (\mu_{i,-1})^2}{2\sigma_i^2}\right) \\ &= \sum_{i=1}^d \underbrace{\left(\frac{\mu_{i,-1} - \mu_{i,+1}}{\sigma_i^2}\right)}_{w_i} x_i + \underbrace{\sum_{i=1}^d \frac{(\mu_{i,+1})^2 - (\mu_{i,-1})^2}{2\sigma_i^2}}_b \\ &= \vec{w}^T \vec{x} + b \end{aligned}$$

$$P(y=+1 | \vec{x}) = \frac{1}{1 + \exp\left(\log\left(\frac{\pi}{1-\pi}\right) + \vec{w}^T \vec{x} + y\right)}$$

$$= \frac{1}{1 + \exp(\vec{w}^T \vec{x} + \underbrace{\log\left(\frac{\pi}{1-\pi}\right) + y}_b)}$$

$$= \frac{1}{1 + \exp(\vec{w}^T \vec{x} + b)} = \frac{1}{1 + e^{(\vec{w}^T \vec{x} + b)}}$$
