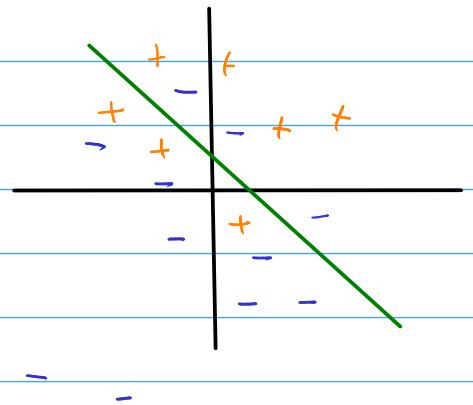
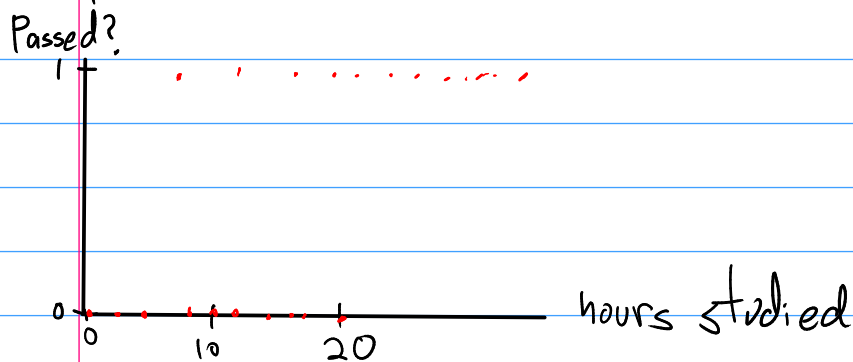


Logistic Regression

Predictions with Confidence



Bernoulli Model

$$\Pr[Y=1|X, w] = p(X, w) \in [0, 1]$$

$$\Pr[Y=0|X, w] = 1 - p(X, w)$$

Take 1: $p(x, w) = \langle x, w \rangle$

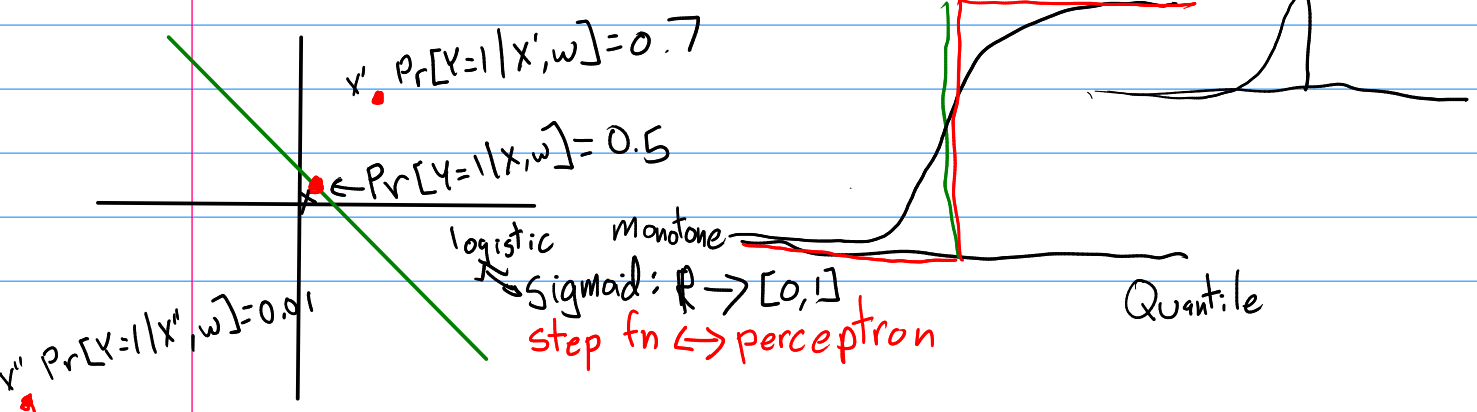
$$\log\left(\frac{p(x, w)}{1 - p(x, w)}\right) = \langle x, w \rangle \Rightarrow \frac{p(x, w)}{1 - p(x, w)} = \exp(\langle x, w \rangle)$$

$$p(x, w) = \exp(\langle x, w \rangle) (1 - p(x, w))$$

$$\Rightarrow p(x, w) (1 + \exp(\langle x, w \rangle)) = \exp(\langle x, w \rangle)$$

$$p(x, w) = \frac{\exp(\langle x, w \rangle)}{1 + \exp(\langle x, w \rangle)} = \frac{1}{1 + \exp(-\langle x, w \rangle)}$$

= sigmoid($\langle x, w \rangle$)



$$p(x, w) = \frac{1}{1 + \exp(-\langle w, x \rangle)}$$

If $p(x, w) > \frac{1}{2}$, predict $\hat{y} = 1$. Else $\hat{y} = 0$.

$$\hat{y} = \text{sign}(\langle w, x \rangle)$$

Maximum Likelihood Estimation

$$\hat{w} = \arg \max_w \prod_{i=1}^n p(x_i, w)^{y_i} (1 - p(x_i, w))^{(1-y_i)}$$

$$p(x_i, w) = \frac{1}{1 + \exp(-\langle w, x_i \rangle)} = p_i$$

$$\rightarrow \arg \max_w \sum \log(p_i^{y_i} (1-p_i)^{(1-y_i)})$$

$$= \arg \max_w \sum y_i \log p_i + (1-y_i) \log(1-p_i)$$

$$\text{Suppose } y_i = 1. \log p_i = \log \left(\frac{1}{1 + \exp(-\langle w, x_i \rangle)} \right) \\ = -\log(1 + \exp(-\langle w, x_i \rangle))$$

$$y_i = 0. \log(1-p_i) = \log \left(\frac{\exp(-\langle w, x_i \rangle)}{1 + \exp(-\langle w, x_i \rangle)} \right) = -\log(1 + \exp(\langle w, x_i \rangle))$$

$$= \arg \max_w \sum \left(-\log(\exp(-y_i \langle w, x_i \rangle) + \exp((1-y_i) \langle w, x_i \rangle)) \right)$$

$$= \arg \min_w \frac{1}{n} \sum \log(\exp(-y_i \langle w, x_i \rangle) + \exp((1-y_i) \langle w, x_i \rangle))$$

$$\tilde{y}_i = \begin{cases} +1 & \text{if } y_i = 1 \\ -1 & \text{if } y_i = 0 \end{cases}$$

$$\rightarrow \arg \min_w \frac{1}{n} \sum \log \left(1 + \exp(-\tilde{y}_i \langle w, x_i \rangle) \right)$$

$$\arg \min_w E_{x, y \sim D} [l_w(x, y)] \approx \arg \min_w \frac{1}{n} \sum l_w(x_i, y_i)$$

$$= \arg \min_w \frac{1}{n} \sum \log(1 + \exp(-\tilde{y}_i \langle w, x_i \rangle))$$

Claim: $l_w(x_i, y_i)$ is convex

Goal: Find \hat{w} where $\frac{1}{n} \sum \nabla_w l(x_i, y_i) = 0$

Claim: $\nabla_w l_w(x_i, y_i) = (p(x_i, w) - y_i)x_i$

No closed form for \hat{w} !

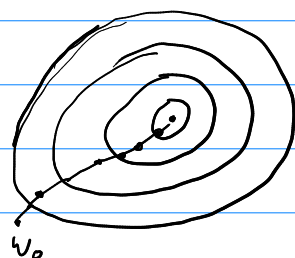
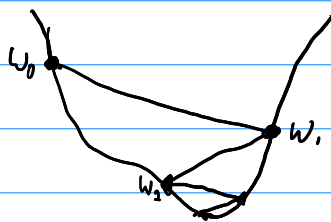
Optimization

Initialize w_0

For $t = 1, 2, \dots$,

Choose: direction d_t
step size η_t

$$w_t = w_{t-1} - \eta_t d_t$$



How to pick: η_t : constant, decaying ($\frac{1}{\sqrt{t}}$), adaptive.

Gradient Descent (GD): $d_t = \frac{1}{n} \sum_{i=1}^n \nabla_w l_{w_{t-1}}(x_i, y_i)$ ($= 0$ at opt).

Stochastic GD (SGD): Pick set $B \subseteq [n]$ randomly

$$d_t = \frac{1}{|B|} \sum_{i \in B} \nabla_w l_{w_{t-1}}(x_i, y_i)$$

Newton's Method: $d_t = \left(\frac{1}{n} \sum \nabla_w^2 l_{w_{t-1}}(x_i, y_i) \right)^{-1} \left(\frac{1}{n} \sum \nabla_w l_{w_{t-1}}(x_i, y_i) \right)$

Inverse Hessian Gradient

Multiclass

$$\Pr[Y=k | X, W] = \frac{\exp(\langle w_k, X \rangle)}{\sum_{k=1}^c \exp(\langle w_k, X \rangle)}$$