

Day -3
June 17, 2024

Last Session

Linear Regression

$$h_{\theta}(x^{(i)}) = \theta^T x^{(i)}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))^2 \quad \text{MSE-Loss}$$

$$\nabla J(\theta) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \rightarrow \text{Grad of MSE-Loss}$$

Gradient Descent

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \eta \cdot \nabla_{\theta_j} J(\theta)$$

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \eta \cdot \left(-\frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \right)$$

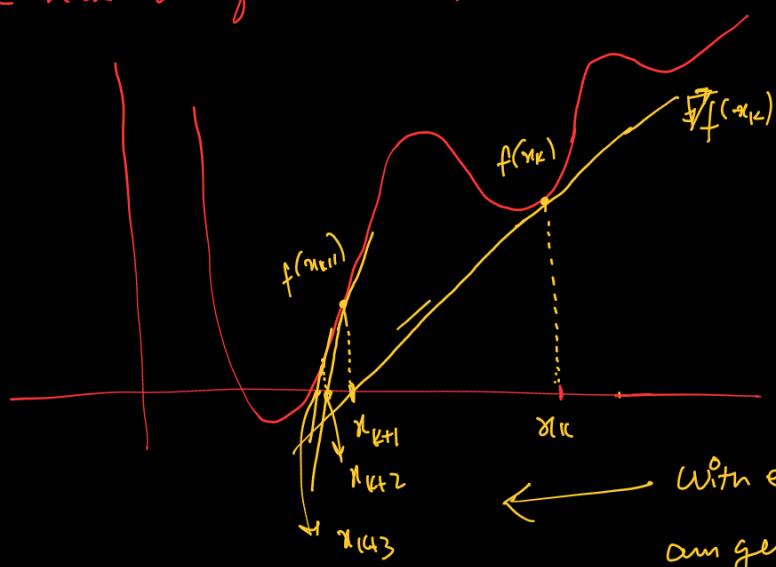
Normal Equations

$$J(\theta) = \frac{1}{m} (Y - X\theta)^T (Y - X\theta)$$

$$\boxed{\theta = (X^T X)^{-1} X^T Y}$$

Newton's Method (Iterative Method) of Optimization.

If we have to find roots of an eqn $f(x) = 0$.



With each moving iteration, I am getting closer & closer to $f(x) = 0$

$$x^{k+1} = x^k - \frac{f(x)}{f'(x)} \rightarrow f(x) = 0 \quad (\text{Iterative process})$$

↳ "Newton's Update Rule"

Optimization

We are not interested in $f(x) = 0$, rather we are interested in finding $\boxed{f'(x) = 0}$

$$\boxed{f(x) = 0 \rightarrow x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}} \quad \boxed{f'(x) = 0 \rightarrow x_{k+1} = x_k - \frac{f'(x_k)}{\boxed{f''(x_k)}} \rightarrow \text{2nd order derivative term}}$$

Newton's Method involves 2nd order derivatives.

→ In general 2nd order methods perform faster than 1st order methods.



Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$.

then given a , we can approximate the pts near " a " (locally) using Taylor's Theorem,

$$\boxed{f(x) \approx f(a) + g^T(x-a) + \frac{1}{2}(x-a)^T H(x-a)}$$

where $g = \nabla f(a)$, $H = \nabla^2 f(a)$

Vector form of Taylor

$$f(x) \approx f(a) + \nabla f(a)^T(x-a) + \frac{1}{2} \underbrace{\nabla^2 f(a)(x-a)^2}_{\geq 0} + \dots$$

← scalar

$$f(x) = f(a) + g^T(x-a) + \frac{1}{2} (x-a)^T H(x-a)$$

$$f(x) = \frac{1}{2} [x^T H x - 2a^T H x + a^T H a] + g^T x - g^T a + f(a)$$

$$f(x) = \frac{x^T H x}{2} - \underbrace{2a^T H x}_{\text{cancel}} + \frac{a^T H a}{2} + \underbrace{g^T x}_{\text{cancel}} - g^T a + f(a)$$

$$= \frac{x^T H x}{2} + g^T x - a^T H x + \frac{a^T H a}{2} - g^T a + f(a)$$

$$= \frac{a^T H a}{2} + (g^T - a^T H) x + \frac{a^T H a}{2} - g^T a + f(a)$$

$$= \boxed{\frac{1}{2} x^T H x + b^T x + c} \rightarrow \text{Want a good eqn in } x.$$

$$\frac{a^T H a}{2} + (g^T - a^T H) x + \frac{a^T H a}{2} - g^T a + f(a)$$

$$\frac{a^T H a}{2} + \overbrace{(g - H^T a)^T}^{\text{cancel}} x + \overbrace{\frac{a^T H a}{2} - g^T a + f(a)}^{\text{cancel}}$$

$$\frac{x^T H a}{2} + \boxed{b^T x} + c$$

$$= b = g - H^T a$$

$$c = \frac{a^T H a}{2} - g^T a + f(a)$$

$$q(x) = \frac{x^T H x}{2} + \underbrace{(g - H^T a)^T x}_{A} + \frac{a^T H a}{2} - g^T a + f(a)$$

$$\nabla q(x) = \cancel{\frac{1}{2} H x} + (g - H^T a) + 0$$

$$\nabla q(x) = H x + (g - H^T a)$$

$$\nabla q(x) = 0$$

$$H x = -(g - H^T a)$$

$$x = H^{-1}(H^T a - g)$$

$$x = H^{-1}(H a - g)$$

(H : Symmetric)

$$\begin{aligned} D(\theta^T H \theta) &\geq 240 \\ \nabla_{\theta} (\theta^T H \theta) &= a \end{aligned}$$

$$x = a - \frac{g}{H}$$

$\boxed{\nabla^2 q(\alpha) = H + \alpha(H - H^T a)}$

$\boxed{\nabla^2 q(\alpha) = H + 0 = H}$

→ "Hessian Matrix"
"Semi-Positive Definite" (PSD)

$\boxed{z^T H z \geq 0}$ Fact

Algorithm

$$\alpha_{t+1} = \alpha_t - \frac{\nabla f(\alpha)}{\nabla^2 f(\alpha)}$$

$$\alpha_{t+1} = \alpha_t - \frac{\nabla J(\theta)}{\nabla^2 J(\theta)}$$

Newton:

$\boxed{\theta^{(t+1)} = \theta^{(t)} - (\nabla^2 J(\theta))^{-1} \nabla J(\theta)}$ Newton's Update rule.

$\downarrow \quad \downarrow$
 $\nabla J(\theta) \quad \nabla^2 J(\theta)$

GD:

$\boxed{\theta^{(t+1)} = \theta^{(t)} - \eta \nabla J(\theta)}$

Probabilistic Interpretation of Linear Regression

In probabilistic interpretation, we are concerned how $y^{(i)}$ are generated given $x^{(i)}$ parameterized by some θ .

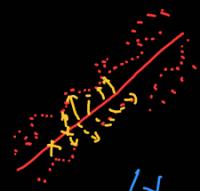
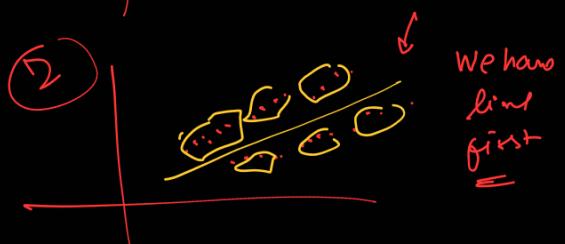
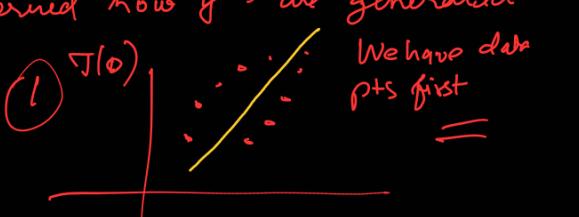
For Linear Regression,

$$y^{(i)} \xrightarrow{\text{approximation}} \theta^T x^{(i)}$$

I am expecting " that my data is distributed around (along) my model $(\theta_0 + \theta_1 x)$ here .

$$y^{(i)} = \theta^T x^{(i)} + \epsilon \rightarrow \epsilon \sim N(0, \sigma^2)$$

↳ "noise"



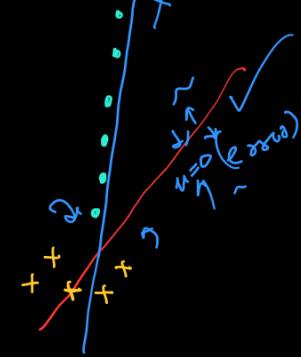
$$y^{(i)} = \theta^T x^{(i)} + \epsilon \quad \epsilon \sim N(0, \sigma^2)$$

Now this becomes the
eqn of how data is generated
given some $x^{(i)}$

$$y^{(i)} = \theta^T x^{(i)} + N(0, \sigma^2)$$

$$(0=0) \quad (\text{mean}=0)$$

This is because we
want on an average the P+
should fall on the line

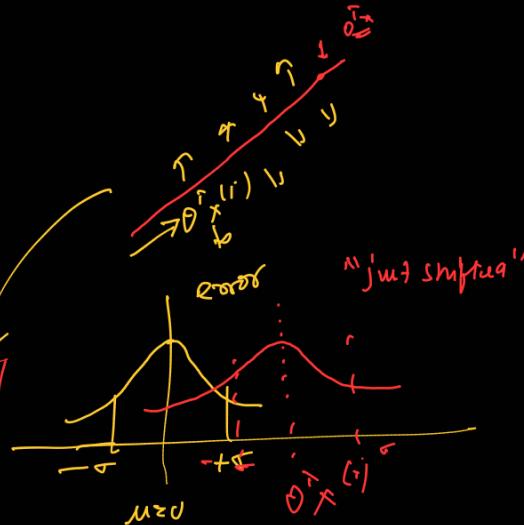


Want on an average the P+
should fall on the line

$$y^{(i)} = \theta^T x^{(i)} + N(0, \sigma^2)$$

This can be interpreted as:-

$$P(y^{(i)} | x^{(i)}; \theta) \sim N(\theta^T x^{(i)}, \sigma^2)$$

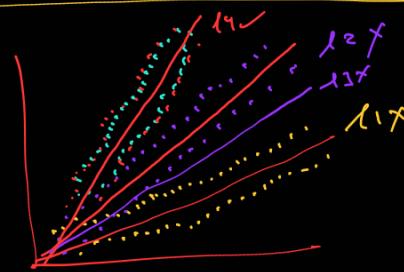
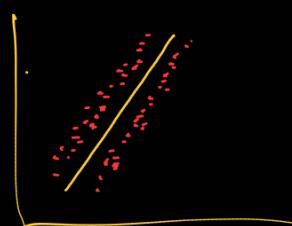


I am assuming that distribution of labels $y^{(i)}$ is Normal, with avg $\theta^T x^{(i)}$,
and some std dev. This also means that on an average, my $y^{(i)}$'s should fall
on the line $\theta^T x^{(i)}$.

$$P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right\} \quad P(z) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(z-\mu)^2}{2\sigma^2}}$$

Now we define the "LIKELIHOOD" of the data given the model.

We are asking "how likely" is this data generated by the given
model.



"We can perform "MLE" to find the model that best generated the data.
(max likelihood estimation)



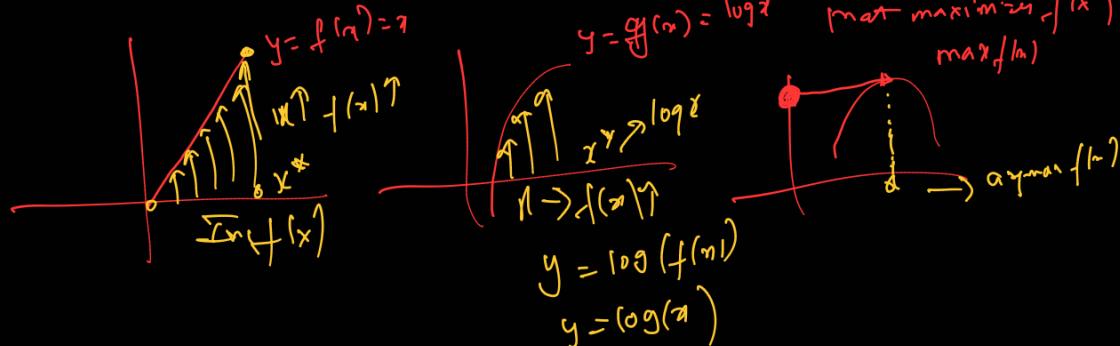
iid

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \underbrace{P(y^{(1)}, y^{(2)}, \dots, y^{(m)})}_{\text{Joint dist}} \underbrace{x^{(1)}, x^{(2)}, \dots, x^{(m)}}_{\text{Joint distribution}},$$

$$f(x) = 2x^2 + 3x + 4$$

$\max f(x)$: I want max value of $f(x)$
 $\operatorname{argmax} f(x)$: I want that "x" that maximizes $f(x)$

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta)$$



$$\hat{\theta}_{MLE} \Rightarrow L(\theta) = \mathcal{L}(\theta)$$

Learn via Tag

"Log likelihood"

$$\hat{\theta}_{MLE} = \mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \log \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta)$$

$$\mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \log \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta)$$

$$\mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}, \theta)$$

$$\mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}} \right)$$

$$\mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \sum_{i=1}^m \log \bar{e}^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$$

$$\mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \bar{e}^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$$

$$\mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$$

$$\frac{f(x)}{4}$$

$$\mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$$

$$L_h(\theta) = \frac{1}{2} \sum_{i=1}^m \theta^T x^{(i)} - \left[\sum_{i=1}^m \left(y^{(i)} - \theta^T x^{(i)} \right)^2 \right]$$

$$L_h(\theta) = \frac{m}{2} \text{argmax}_{\theta} \left[\sum_{i=1}^m \left(y^{(i)} - \theta^T x^{(i)} \right)^2 \right]$$

$$L_h(\theta) = \frac{m}{2} \text{argmax}_{\theta} - J(\theta)$$

$$\min f(x) \underset{\text{max}}{=} -f(x)$$

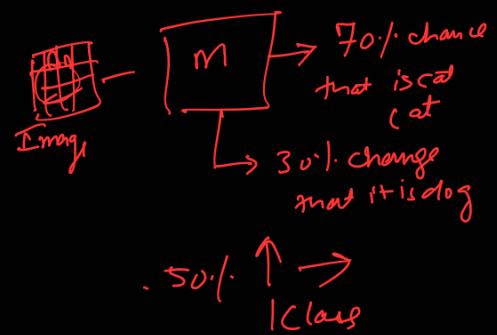
$$\max f(x) \underset{\text{min}}{=} -f(x)$$

$$\boxed{L_h(\theta) \underset{\text{max}}{\uparrow} -J(\theta) \underset{\text{max}}{\uparrow} J(\theta) \underset{\text{minimize}}{\downarrow}}$$

↳ We were right about Probabilistic Int of Linear Regression
these both are "interchangeable"

Logistic Regression (Classification) \rightarrow (Cat vs Dog)

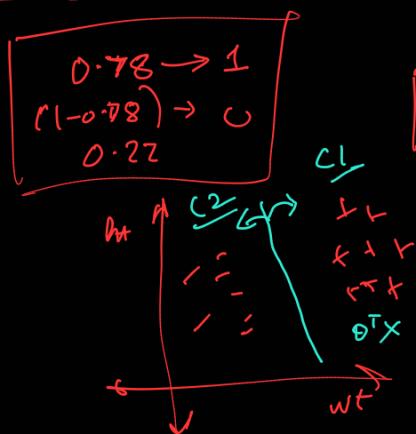
Linear Regression \rightarrow continuous domain:
100.28, 75.67, 101.82, 224 ...



Logistic Regression \rightarrow Labels depending on how many categories I have.

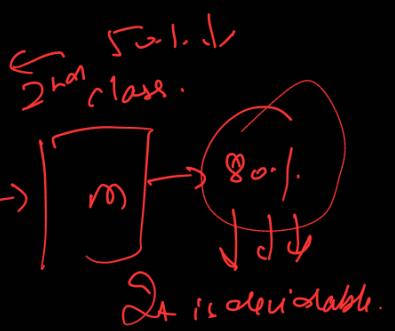
Binary Classification (2 Categories)

$$y \in \{0, 1\}$$



$$\boxed{y^{(i)} = \{0, 1\}} \rightarrow (p)(1-p)$$

then $y^{(i)} | x^{(i)}, \theta \sim \text{Bernoulli}(p)$



$$\text{Bernoulli}(\phi) = 1 \text{ if } \phi > 0.5$$

$$\text{Multinomial}(\phi_m) = 1$$

$\{0\} \leftarrow \text{fp}$



Binomial distribution
↳ Repeated Bernoulli trials

$$\propto n! p^n q^{n-n}$$

$$= n! n! p^n (1-p)^{n-n}$$

Multinomial distribution

↳ Repeated multinomial trials

8 Teams

$$\frac{T_1! T_2! \dots T_8!}{(T_1 + T_2 + \dots + T_8)!}$$

ϕ_m (0,0) modulated

Let $d^{(i)}$ be the distance of the line and the point $(x^{(i)})$, then ideally to have a high probability for a point that is further from the line (E.g. - ureq.)

$$-\infty < d^{(i)} < \infty \xrightarrow{\text{map}} g(d) \in [0, 1]$$

Actually I can mapping my distances from the line into probabilities.

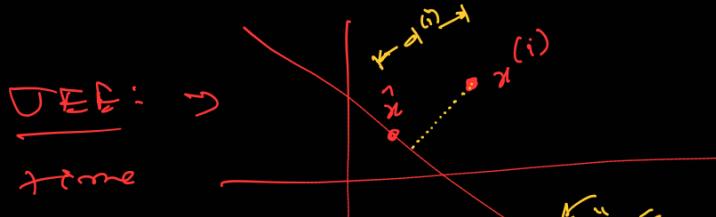
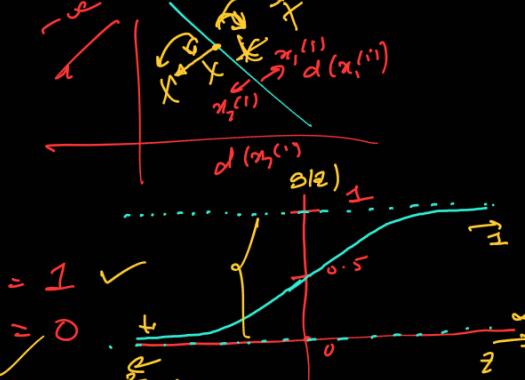
$$d^{(i)} = \omega^T (\mathbf{x}^{(i)} - \mathbf{z})$$

$$d^{(i)} = \omega^T \mathbf{x}^{(i)} - \omega^T \mathbf{z} \quad \xrightarrow{\text{constant}}$$

$$d^{(i)} = \omega^T \mathbf{x}^{(i)} + b$$

differentiable ✓
 $(0, 1)$ ✓
 $\mathbf{z} \rightarrow \infty \rightarrow g(z) = 1$ ✓
 $\mathbf{z} \rightarrow -\infty \rightarrow g(z) = 0$ ✓

$$\text{Sigmoid: } \sigma(z) = \frac{1}{1 + e^{-z}}$$



$$d^{(i)} = \omega^T (\mathbf{x}^{(i)} - \hat{\mathbf{x}})$$

$$\omega^T \mathbf{x}^{(i)} + b = \underline{\omega^T \mathbf{x}^{(i)} + b}$$

$$P(y^{(i)} = 1 | \mathbf{x}^{(i)}, \theta) = g(d^{(i)})$$

$$\text{Fact: } g(z) = \frac{1}{1 + e^{-z}}$$

$$P(y^{(i)} = 1 | \mathbf{x}^{(i)}, \theta) = g(\omega^T \mathbf{x}^{(i)})$$

$$g(z) = \sigma(z) (1 - \sigma(z))$$

$$\therefore g(z) = \frac{1}{1 + e^{-z}}$$

$$\sigma(-z) = 1 - \sigma(z)$$

$$1 - \frac{1}{1 + e^{-z}} = \frac{1 + e^z - 1}{1 + e^z}$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\sigma(z) = (1 + e^{-z})^{-1}$$

$$\sigma'(z) = -\cancel{x} \cdot (1 + e^{-z})^{-2} \cdot (0 + e^{-z}) \cdot (-1)$$

$$= \frac{e^{-z}}{(1 + e^{-z})^2}$$

$$= \frac{e^{-z}}{(1 + e^{-z})} \cdot \frac{1}{(1 + e^{-z})}$$

$$\left(\frac{1}{1 + e^{-z}}\right) \cdot \left(\frac{1}{1 + e^{-z}}\right)$$

$$= \left(\frac{1 + e^{-z} - e^{-z}}{1 + e^{-z}}\right) \cdot \left(\frac{1}{1 + e^{-z}}\right)$$

$$= \left(\left(\frac{1 + e^{-z}}{1 + e^{-z}}\right) - \frac{e^{-z}}{1 + e^{-z}}\right) \cdot \left(\frac{1}{1 + e^{-z}}\right)$$

$$= \left(1 - \frac{1}{1 + e^{-z}}\right) \left(\frac{1}{1 + e^{-z}}\right)$$

$$\boxed{(1 - \sigma(z)) \times \sigma'(z)}$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

$$\boxed{\begin{aligned} P(y^{(i)} = 1 | \mathbf{x}^{(i)}, \theta) &= g(\mathbf{x}^{(i)} \cdot \theta) \\ &= \frac{1}{1 + e^{-\theta^T \mathbf{x}^{(i)}}} \quad \checkmark \end{aligned}}$$

$$\boxed{\begin{aligned} P(y^{(i)} = 0 | \mathbf{x}^{(i)}, \theta) &= 1 - P(y^{(i)} = 1 | \mathbf{x}^{(i)}, \theta) \\ &= 1 - \frac{1}{1 + e^{-\theta^T \mathbf{x}^{(i)}}} \quad \checkmark \end{aligned}}$$

$$\text{Bernoulli} \stackrel{?}{=} P(y = 1) \cdot P(y = 0) = \frac{1}{2} \cdot \frac{1}{2} \xrightarrow{\text{Fact 1}} \underline{\underline{\text{Fact 1}}}$$

Log likelihood

$$L = \prod_{i=1}^m \phi^{y(i)} (1-\phi)^{1-y(i)}$$

$$\ln L(\phi) = \sum_{i=1}^m y(i) \ln \phi + (1-y(i)) \ln (1-\phi)$$

$P(y^{(i)} = 1 | \alpha^{(i)}, \theta)$ $P(y^{(i)} = 0 | \alpha^{(i)}, \theta)$

$$L(\theta) = \prod_{i=1}^m P(y^{(i)} = 1 | \alpha^{(i)}, \theta) P(y^{(i)} = 0 | \alpha^{(i)}, \theta)$$

$$L(\theta) = \prod_{i=1}^m g(\theta^T \alpha^{(i)}) \cdot (1 - g(\theta^T \alpha^{(i)}))$$

(AND PESUM)

$$L(\theta) = \prod_{i=1}^m h_\theta(\theta^T \alpha^{(i)}) (1 - h_\theta(\theta^T \alpha^{(i)}))^{(1-y(i))}$$

$$\ln L(\theta) = \log \prod_{i=1}^m h_\theta(\theta^T \alpha^{(i)}) (1 - h_\theta(\theta^T \alpha^{(i)}))^{(1-y(i))}$$

$$\ln L(\theta) = \sum_{i=1}^m \log h_\theta(\theta^T \alpha^{(i)}) (1 - h_\theta(\theta^T \alpha^{(i)}))^{(1-y(i))}$$

$$\ln L(\theta) = \sum_{i=1}^m \log \tau(\theta^T \alpha^{(i)}) (1 - \tau(\theta^T \alpha^{(i)}))^{(1-y(i))}$$

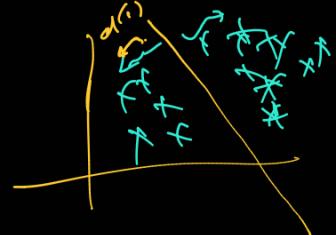
$$\nabla_{\theta} \ln L(\theta) = \sum_{i=1}^m \frac{y(i) / \log \tau(\theta^T \alpha^{(i)}) + (1-y(i)) \log (1 - \tau(\theta^T \alpha^{(i)}))}{\tau(\theta^T \alpha^{(i)}) (1 - \tau(\theta^T \alpha^{(i)}))}$$

$g(f(x)) = g'(f(x)) \cdot f'(x)$

$$\nabla_{\theta} \ln L(\theta) = \sum_{i=1}^m \left[y(i) \left(\frac{1}{\tau(\theta^T \alpha^{(i)})} - \frac{(1-y(i))}{1-\tau(\theta^T \alpha^{(i)})} \right) \tau'(\theta^T \alpha^{(i)}) (1 - \tau(\theta^T \alpha^{(i)})) \right] \alpha^{(i)}$$

$$\nabla_{\theta} \ln L(\theta) = \sum_{i=1}^m \left[y(i) \left(\frac{1}{\tau(\theta^T \alpha^{(i)})} + \frac{(1-y(i))}{1-\tau(\theta^T \alpha^{(i)})} \right) \tau'(\theta^T \alpha^{(i)}) \right] \alpha^{(i)}$$

$$\nabla_{\theta} \ln L(\theta) = \sum_{i=1}^m \left[y(i) \left(\frac{1}{\tau(\theta^T \alpha^{(i)})} + \frac{(1-y(i))}{1-\tau(\theta^T \alpha^{(i)})} - \tau'(\theta^T \alpha^{(i)}) \right) \right]$$



$$\mathcal{J}_{\theta} L_{\text{L}}(\theta) = \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)}) \sigma^{(i)}$$

$$\mathcal{J}_{\theta} L_h(\theta) = \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) \sigma^{(i)}$$

MSE loss derivative = $\left[\frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) \sigma^{(i)} \right]$

logistic Regression and Linear Regression have similar properties.

$$D_j^{(t+1)} \leftarrow D_j^{(t)} - \eta \sigma_{\theta} \nabla_{\theta} J(\theta)$$

$$h_{\theta}(x^{(i)}) = \frac{(y^{(i)} - \theta^T x^{(i)})}{1 + e^{-\theta^T x^{(i)}}}$$

Generalized Class of Linear Models (GCLM)

$$L(\theta) = \prod_{i=1}^m P(y^{(i)} = 1 | x^{(i)}, \theta)^{y^{(i)}} P(y^{(i)} = 0 | x^{(i)}, \theta)^{1-y^{(i)}}$$

$$LL(\theta) = \log \prod_{i=1}^m P(y^{(i)} = 1 | x^{(i)}, \theta)^{y^{(i)}} P(y^{(i)} = 0 | x^{(i)}, \theta)^{1-y^{(i)}}$$

$$LL(\theta) = \sum_{i=1}^m \log P(y^{(i)} = 1 | x^{(i)}, \theta)^{y^{(i)}} + \sum_{i=1}^m \log P(y^{(i)} = 0 | x^{(i)}, \theta)^{1-y^{(i)}}$$

$$\sum_{i=1}^m y^{(i)} \log \underbrace{P(y^{(i)} = 1 | x^{(i)}, \theta)}_{\text{HYPOTHESIS}} + (1-y^{(i)}) \log \underbrace{P(y^{(i)} = 0 | x^{(i)}, \theta)}_{g(h_{\theta}(x^{(i)}))}$$

$$\begin{cases} h_{\theta}(x^{(i)}) = \theta^T x^{(i)} \\ P(y^{(i)} = 1 | x^{(i)}, \theta) = g(h_{\theta}(x^{(i)})) = g(\theta^T x^{(i)}) \end{cases}$$

$$\sum_{i=1}^m y^{(i)} \rho(\theta^T x^{(i)}) + (1-y^{(i)}) \log [1 - \sigma(\theta^T x^{(i)})]$$

$$\begin{aligned}
&= \sum_{i=1}^m y^{(i)} \log \sigma(\widehat{\theta}_{\alpha}^{(i)}) + (1-y^{(i)}) \log (1-\sigma(\widehat{\theta}_{\alpha}^{(i)})) \\
L(\theta_0) &\geq \sum_{i=1}^m y^{(i)} \log \sigma(\widehat{\theta}_{\alpha}^{(i)}) + (1-y^{(i)}) \log (1-\sigma(\widehat{\theta}_{\alpha}^{(i)})) \\
\partial_{\theta} L(\theta) &= \sum_{i=1}^m \frac{y^{(i)}}{\sigma(\widehat{\theta}_{\alpha}^{(i)})} \left(\frac{(1-y^{(i)})}{1-\sigma(\widehat{\theta}_{\alpha}^{(i)})} \cdot x^{(i)} \right) + \frac{(1-y^{(i)})}{1-\sigma(\widehat{\theta}_{\alpha}^{(i)})} \cdot \left(\frac{(1-y^{(i)})}{1-\sigma(\widehat{\theta}_{\alpha}^{(i)})} \cdot x^{(i)} \right)^T \\
\partial_{\theta} h(\theta) &= \sum_{i=1}^m y^{(i)} \left(\frac{(1-y^{(i)})}{1-\sigma(\widehat{\theta}_{\alpha}^{(i)})} \cdot x^{(i)} \right) + \frac{(1-y^{(i)})}{(1-\sigma(\widehat{\theta}_{\alpha}^{(i)}))} \left(-\sigma'(\widehat{\theta}_{\alpha}^{(i)}) / (1-\sigma(\widehat{\theta}_{\alpha}^{(i)})) \right) x^{(i)} \\
&= \sum_{i=1}^m x^{(i)} \left(y^{(i)} - \sigma(\widehat{\theta}_{\alpha}^{(i)}) \right) + \frac{y^{(i)}}{\sigma(\widehat{\theta}_{\alpha}^{(i)})} \left(\sigma'(\widehat{\theta}_{\alpha}^{(i)}) - \sigma(\widehat{\theta}_{\alpha}^{(i)}) \right) \\
&= \sum_{i=1}^m x^{(i)} \left(y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}) \right) \quad \xrightarrow{\text{Logistic Loss}} \text{Logistic Loss} \\
&= - \sum_{i=1}^m x^{(i)} \left(y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}) \right) \quad \text{on } \mathcal{S} \\
&\quad \text{Identity } |\mathcal{S}| = n \\
&\quad \text{---} \quad \text{---} \quad \text{---} \\
&\quad \boxed{S(f(a)) \quad g(f(a))} \\
&\quad \boxed{S = n}
\end{aligned}$$