

L4

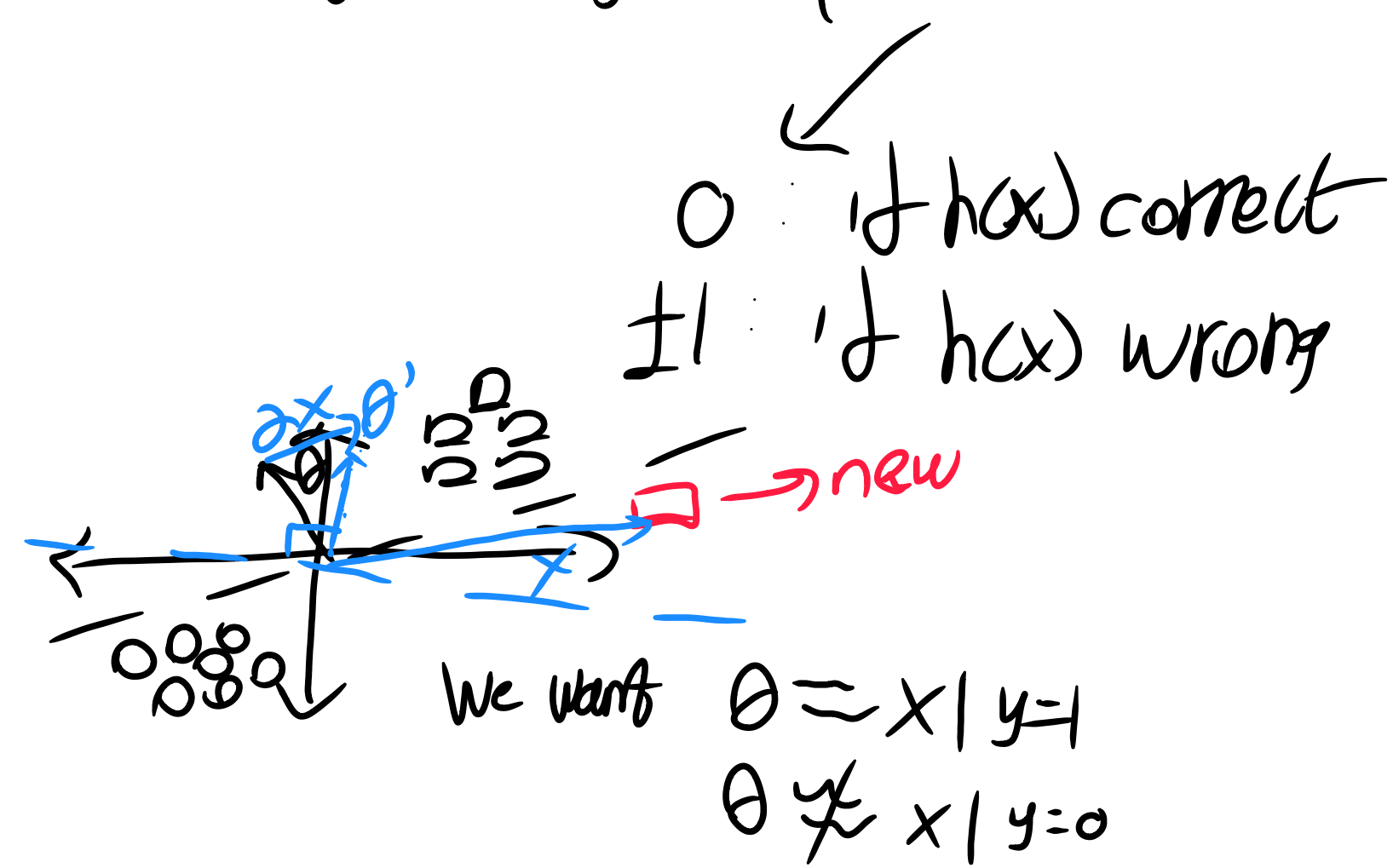
Perceptron algorithm

Similar to sigmoid

$$g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

$h(x) = g(\theta^T x)$

$\theta_j := \theta_j + \alpha(y^{(i)} - h(x^{(i)}))x_j^{(i)}$



perceptron algorithm: add data one by one

disadvantage: $\frac{x}{0}$ never classify as no boundary

Exponential family

$p(y; \eta) = b(y) e^{\eta^T T(y) - a(\eta)} = b(y) \frac{e^{\eta^T T(y)}}{e^{a(\eta)}}$

y data η natural parameter
 $T(y)$ sufficient statistic
 $b(y)$ - base measure
 $a(\eta)$ - log-partition

} match dimension

eg. Bernoulli: $p(y; \phi) = \phi^y (1-\phi)^{1-y}$
 $= e^{y \log(\phi) + (1-y) \log(1-\phi)}$
 $= e^{y \log(\phi) - y \log(\phi) + \log(1-\phi)}$

Gaussian with $\sigma^2 = 1$

$p(y; \mu) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2}}$
 similar

Why exponential family?

① MLE wrt $\eta \Rightarrow$ concave
 NLL is convex

② $E[y; \eta] = \frac{\partial a(\eta)}{\partial \eta}$

③ $\text{Var}(y; \eta) = \frac{\partial^2 a(\eta)}{\partial \eta^2}$

If η vector $\Rightarrow \text{Var} = H(a(\eta))$

GLM

Assumptions

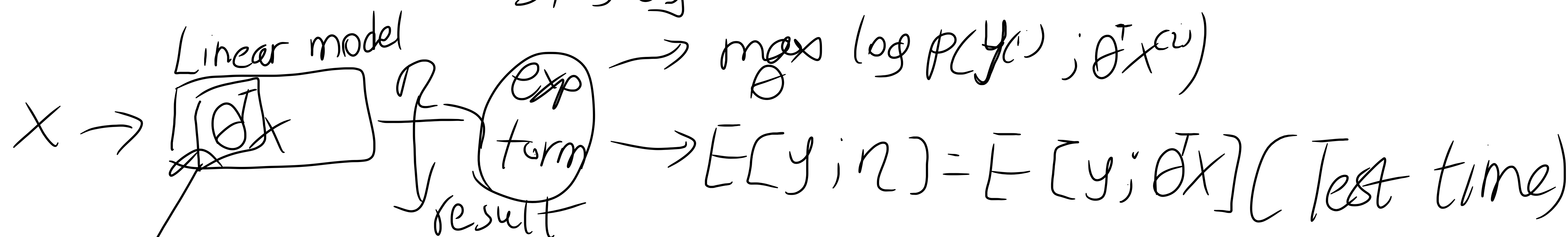
1. $y|x, \theta \sim$ exponential family

Real - Gaussian
 Binary - Bernoulli
 Count - Poisson
 \mathbb{R}^+ - Gamma, Exponential
 Dirich - Beta, Dirichlet } Bayesian

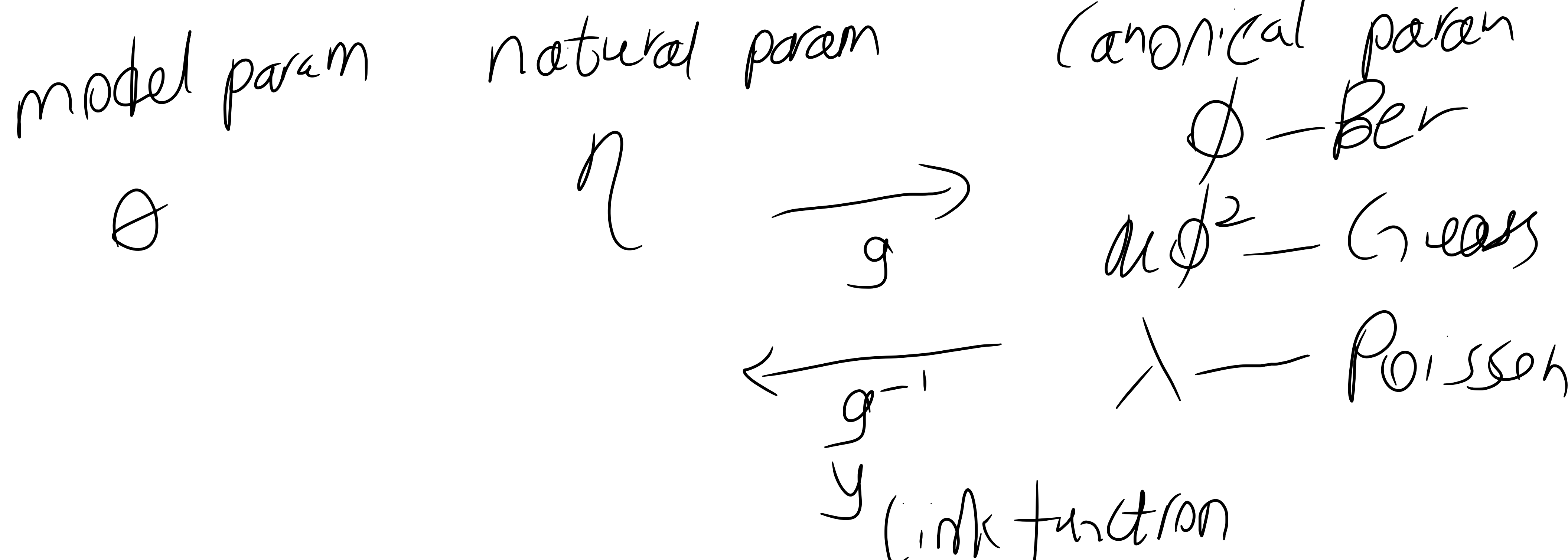
2. $\eta = \theta^T x, \theta \in \mathbb{R}^n, x \in \mathbb{R}^n$

3. Test time output $E[y|x; \theta]$

$\Rightarrow h_\theta(x) = E[y|x; \theta]$



3 - Parameter



GLM: Learning update rule same for exp family

$\theta_j := \theta_j + \alpha(y^{(i)} - h_\theta(x^{(i)}))x_j^{(i)}$

$E[y; \eta] = g(\eta) \rightarrow$ canonical response function

$\mu = E[y; \eta] = g(\eta)$
 $g(\eta) = \frac{\partial a(\eta)}{\partial \eta}$

$\Rightarrow g(\eta) = \eta$

Logistic regression

$h_\theta(x) = E[y; \eta] = \phi = \frac{1}{1+e^{-\eta}}$

Softmax Regression

cross-entropy minimization

$k = \# \text{ classes}$

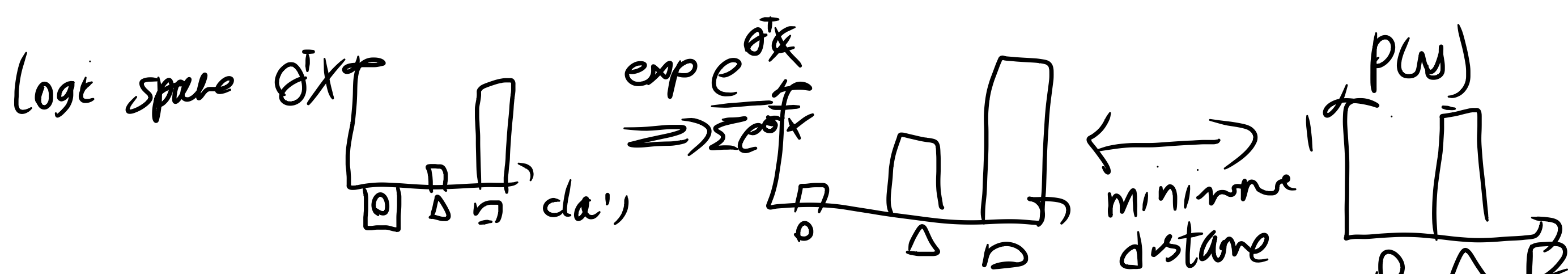
Labels $y = \{1, \dots, k\}$

θ_{Δ}

$\theta_{\Delta} \in \mathbb{R}^n$

$k \neq$ such that

θ_{Δ} class $\in \{\Delta, 0, \square, \dots\}$



Cross Entropy $(P, \hat{P}) = -\sum p(y) \log \hat{p}(y)$

minimize Gradient descent