# Random features and kernels

Goals: law of large numbers
random features as kernels

Office hr tomorrow   11 - noon

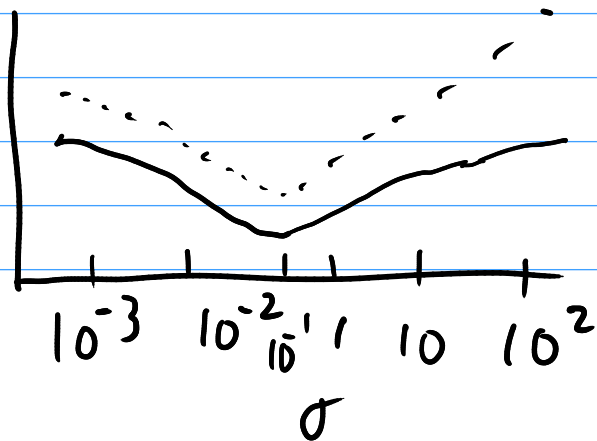# Homework Q's?

**1.4**
$$\nabla_w L = \frac{-1}{n} \sum_{i=1}^{n} P\left[ Y_i = -y_i \mid \vec{x}_i \right] \vec{x}_i y_i$$

$\underbrace{\qquad\qquad\qquad}$ need practical form

**2/3**
$$\min_{\vec{w}, b} \; \boxed{L_{logistic}(\vec{w}, b)} + \lambda \|\vec{w}\|^2$$

cost = objective function

**3**



after convergence

$10^{-3} \quad 10^{-2} \; 10^{-1} \; 1 \quad 10 \quad 10^2$

$\sigma$

$$\begin{pmatrix} \vec{w} \\ b \end{pmatrix}$$
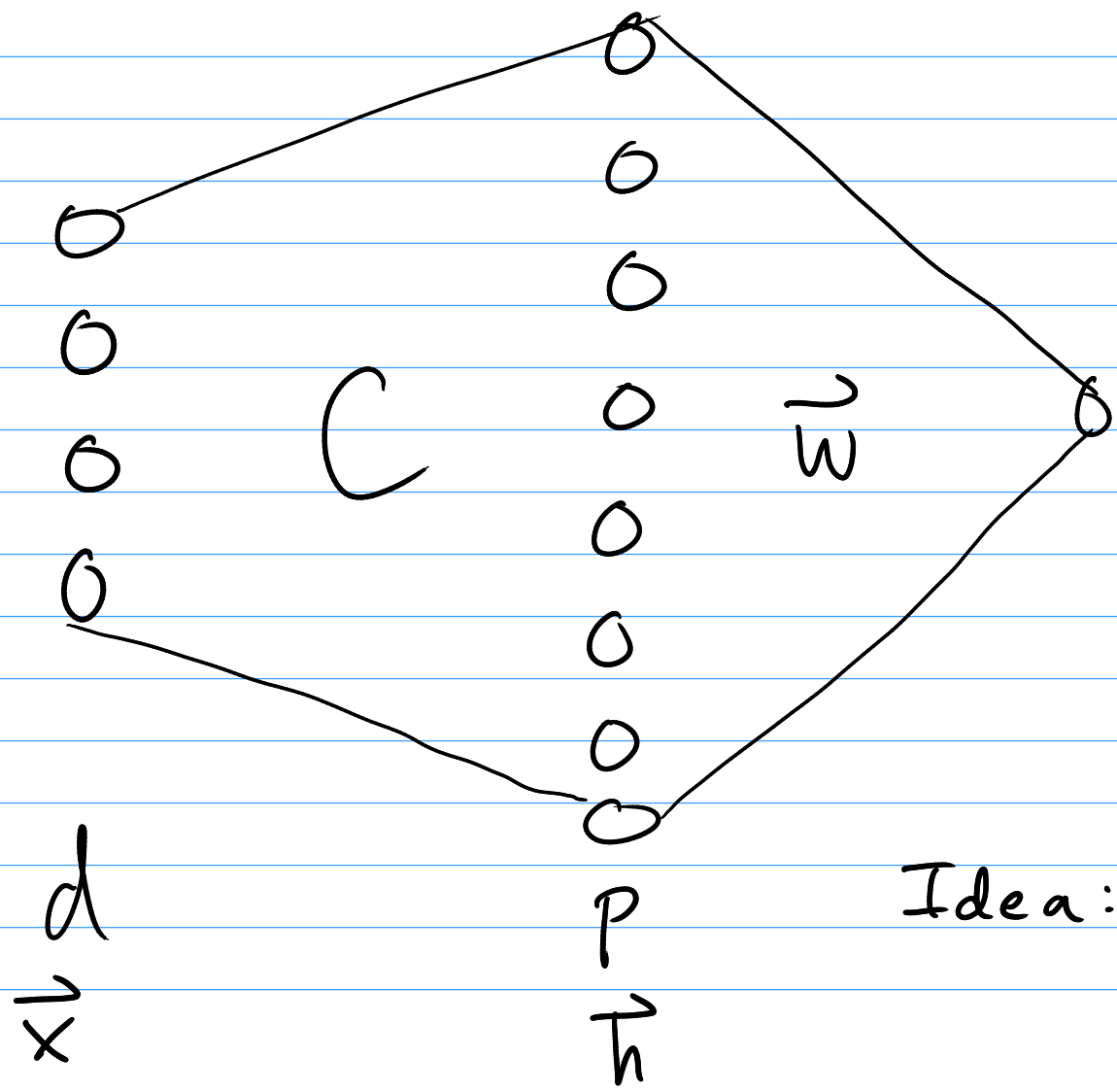
loop

   evaluate gradient at $\vec{w}_t, b_t$

   $\vec{w}_{t+1} \leftarrow \vec{w}_t + h \cdot \text{grad } w$

   $b_{t+1} \leftarrow b_t + h \cdot \text{grad } b$

   $\| \vec{w}_{t+1} - \vec{w}_t \|_\infty$

# Random Features

Rahimi
Recht 2008
Neal, Williams
1990's

2-layer NN



$d$

$\vec{x}$

$p$

$\vec{h}$

$$\vec{h} = g(C\vec{x})$$

$$h_i = g(\vec{c_i}^T \vec{x})$$

↑ $i$-th row

$$f(\vec{x}) = \sum_{i=1}^{p} w_i h_i$$

Idea: Think about $p \to \infty$
wide limit

Fix $C$ to its random init

Going to study geometry of hidden layer representations

inner products

Take $\vec{x}, \vec{x}' \in \mathbb{R}^d$ two inputs

Study $\frac{1}{P} \underbrace{\hat{h}(\vec{x})^T \hat{h}(\vec{x}')}_{\in \mathbb{R}^P} = \frac{1}{P} \sum_{i=1}^{P} h_i(\vec{x}) h_i(\vec{x}')$

$$= \frac{1}{P} \sum_{i=1}^{P} g(\vec{c_i}^T \vec{x}) \, g(\vec{c_i}^T \vec{x}')$$

random

Want to
take $P \to \infty$ ... renormalize sum
think of

$$h_i \longrightarrow \frac{1}{\sqrt{P}} h_i$$

like rescaling $c$'s to be $O\left(\frac{1}{\sqrt{P}}\right)$

# Law of Large Numbers (weak law)

Collection of random vars $\in \mathbb{R}$ $\qquad X_i$ : iid. $\quad i = 1, \ldots, P$

$$Var[X_i] = \sigma^2 < \infty$$
$$\mathbb{E}[(X_i - \mathbb{E}[X_i])^2]$$

Partial sum

random var

$$S_p = \frac{1}{P} \sum_{i=1}^{P} X_i \xrightarrow[P \to \infty]{\text{in probability}} \mathbb{E}[X_i] \quad \begin{array}{c} \text{not} \\ \text{random} \end{array}$$

$$\hat{S}_p = \frac{1}{P} \sum_{i=1}^{P} s_i \qquad \qquad \overset{``}{\int} X_i \, dP(X_i)$$

$$Pr\left[ |S_p - \mathbb{E}[X_i]| > \varepsilon \right] \searrow 0 \quad \text{as } p \to \infty$$
$$\text{for all } \varepsilon > 0$$

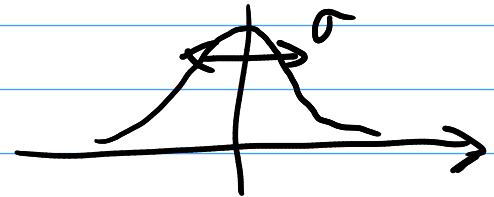Concentration inequality (ML, probability useful)
  Markov, Hoeffding, etc.

Chebyshev's inequality  $X$ r.v.  $\text{Var}[X] < \infty$

$$\Pr\left[|X - \mathbb{E}[X]| \geq a\right] \leq \frac{\text{Var}[X]}{a^2}$$

ex/ if $a = 5 \cdot \sqrt{\text{Var}[X]}$

r.h.s. $= \frac{\text{Var}[X]}{\left(5 \cdot \sqrt{\text{Var}[X]}\right)^2} = \frac{1}{25}$

Use Chebyshev on $S_P = \frac{1}{P} \sum_{i=1}^{P} X_i$

Compute $\mathbb{E}[S_P] = \mathbb{E}\left[\frac{1}{P} \sum_{i=1}^{P} X_i\right] = \frac{1}{P} \sum_{r=1}^{P} \mathbb{E}[X_r] = \mathbb{E}[X_r]$

$$Var[S_P] = Var\left[\sum_{i=1}^{P} \left(\frac{X_i}{P}\right)\right]$$

$$= \sum_{i=1}^{P} Var\left[\left(\frac{X_i}{P}\right)\right]$$

$$= \sum_{i=1}^{P} \left(\frac{1}{P}\right)^2 \underbrace{Var[X_i]}_{\sigma^2}$$

$$= P \cdot \left(\frac{1}{P}\right)^2 \cdot \sigma^2$$

$$= \frac{\sigma^2}{P}$$

$Var[A+B]$
$= Var[A] + Var[B]$
  if $A, B$ indep.

$Var[zA] = z^2 Var[A]$

Variance of Sample averages decrease as $\frac{1}{P}$

Using Chebyshev,

$$Pr\left[\,|S_p - \mathbb{E}[X_i]|\,\geq\,\varepsilon\,\right]\,\leq\,\frac{\sigma^2}{p\,\varepsilon^2}\,\searrow\,0$$

Ex/ 99% confidence $\Rightarrow$ Pr [further than $\varepsilon$] $<$ 1%

$$\frac{\sigma^2}{p\,\varepsilon^2}\,<\,0.01\,\Rightarrow\,\varepsilon\,>\,\frac{\sigma}{\sqrt{p}}\,\frac{1}{\sqrt{0.01}}$$

tolerance $\sim\,\frac{1}{\sqrt{p}}$

Returning to random features

$$\frac{1}{p} \vec{h}(\vec{x})^T \vec{h}(\vec{x}') = \frac{1}{p} \sum_{i=1}^{p} g(\vec{c_i}^T \vec{x}) g(\vec{c_i}^T \vec{x}')$$

iid $\vec{c_i}$

$$\xrightarrow{P} \mathbb{E}_{\vec{c}} \left[ g(\vec{c}^T \vec{x}) g(\vec{c}^T \vec{x}') \right]$$

$$:= K(\vec{x}, \vec{x}')$$ random feature kernel

geometry of network

kernel like in ridge regression/SVM

works for $g$ not too crazy (smooth)

$$\vec{c_i} \overset{iid}{\sim} \mu$$

$$\mathbb{E}\left[ \|\vec{c}\|^2 \right] < \infty$$