

$$\Pr(X=k) = \binom{m}{k} p^k (1-p)^{m-k}$$

$p = \Pr(X \geq k)$ = The prob under the null hypo that we would see at least this many flies?

$$= \sum_{i=k}^m \binom{m}{k} \frac{1}{2}^k \frac{1}{2}^{m-k} = \frac{1}{2}^m \sum_{i=k}^m \binom{m}{k}$$

Ex. - $m=10, k=9, p=0.01$	$m=10, k=5, p=0.62$
- $m=100, k=90, p \approx 10^{-6}$	
- $m=100, k=50, p=0.54$	

Reject H_0 if evidence unlikely under H_0 distribution

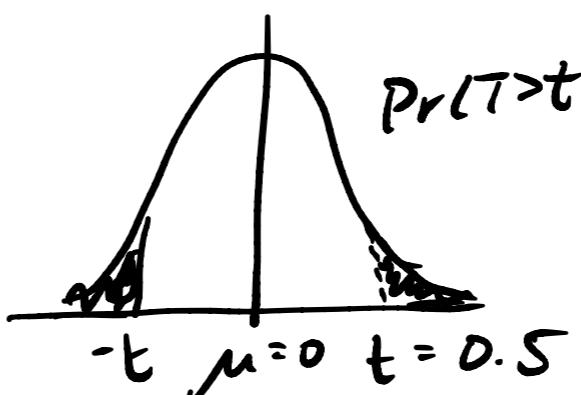
		D_1	D_2	...	D_m	$d_i = \text{error}_2 - \text{error}_1$
		a_1	0.1	0.05	...	0.05
		a_2	0.2	0.03	...	0.07
d_i		+0.1	-0.02	-	-	+0.02
T is our RV,		= normalized average.				
		differences in error.				

$$\bar{d} = \frac{1}{m} \sum_{i=1}^m d_i \quad S_d = \sqrt{\frac{1}{m-1} \sum_{i=1}^m (d_i - \bar{d})^2}$$

$$t = \frac{\bar{d} - \mu}{S_d / \sqrt{m}}$$

Assume T follows a Student t-distribution
 μ = expect average diff.

$$H_0: \mu = 0 \quad H_1: \mu > 0 \quad (\text{as } \mu < \mu \text{ is better})$$



$p = \Pr(T > t)$ if p is small, the evidence is unlikely under H_0 , so reject H_0 .

$p = \Pr(T > t) \text{ or } \Pr(T < -t)$

Nov. 26, 2019: find (latent) features $h \in \mathbb{R}^k$ for $x \in \mathbb{R}^d$
that explain x and have certain properties

$k < d \rightarrow$ goal: lower-dimensional representation
compactness, avoiding overfitting

$$(1) \quad \begin{aligned} &\text{Learn dictionary } D \in \mathbb{R}^{k \times d} \quad x \approx hD = \vec{h}_1 \vec{d}_1 + \vec{h}_2 \vec{d}_2 + \dots + \vec{h}_k \vec{d}_k \\ &\min_{D, h_1, \dots, h_n} \sum_{i=1}^n \|x_i - h_i D\|_2^2 \quad \begin{aligned} &\text{Found } D \quad \xrightarrow{\text{scalars}} \\ &\text{and } h_1, \dots, h_n \quad \phi(x_i) = h_i \end{aligned} \end{aligned}$$

\vec{h}_i
vectors

For new x , $\phi(x) = \arg \min_h \|x - hD\|_2^2$

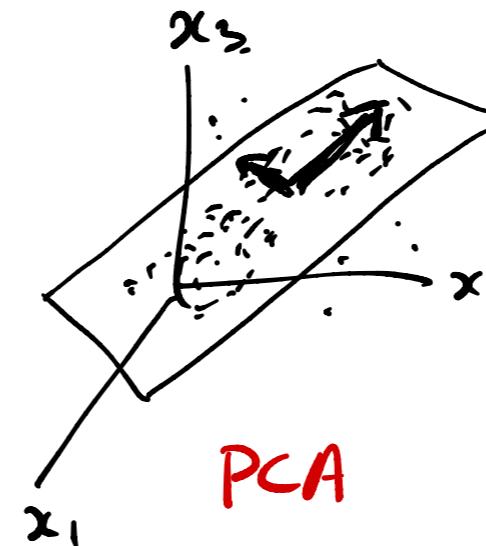
Solution to (1): Let $X = U\Sigma V^T \in \mathbb{R}^{n \times d}$

$$D = V_{:, 1:k}^T \quad k \times d$$

$$h = \begin{bmatrix} h_1 \\ \vdots \\ h_n \end{bmatrix} = U_{:, 1:k} \Sigma_k$$

$$\underline{h_i = U_{i, 1:k} \Sigma_k}$$

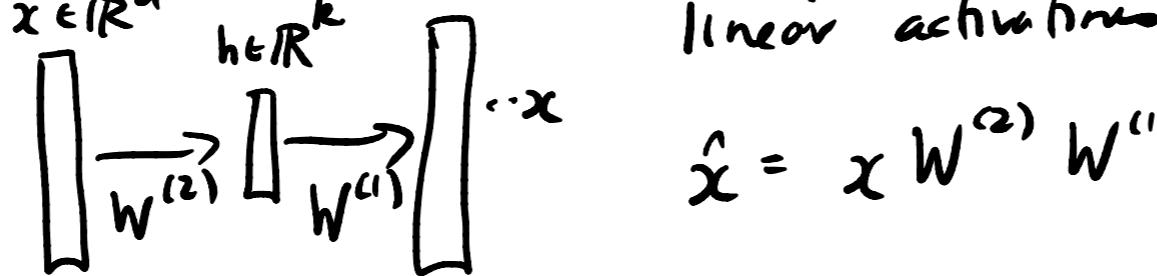
$$\Sigma = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_d \end{bmatrix} \quad \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d \geq 0$$



How is this different from NNs? Equivalent for one setting

linear autoencoder

produce $\hat{x} \approx x$.



linear activation

$$\hat{x} = x W^{(2)} W^{(1)}$$

$$\min_{W^{(1)}, W^{(2)}} \sum_{i=1}^n \| (\underbrace{x_i \cdot W^{(2)}}_{\hat{x}_i}) W^{(1)} - x_i \|_2^2 \rightarrow W^{(2)} = V_{:, 1:k}$$

$$W^{(1)} = V_{:, 1:k}^T$$

$$h_i := x_i \cdot W^{(2)} = U_{i,:} \sum V^T V_{:, 1:k} = \underline{U_{i,:} \sum_k}$$

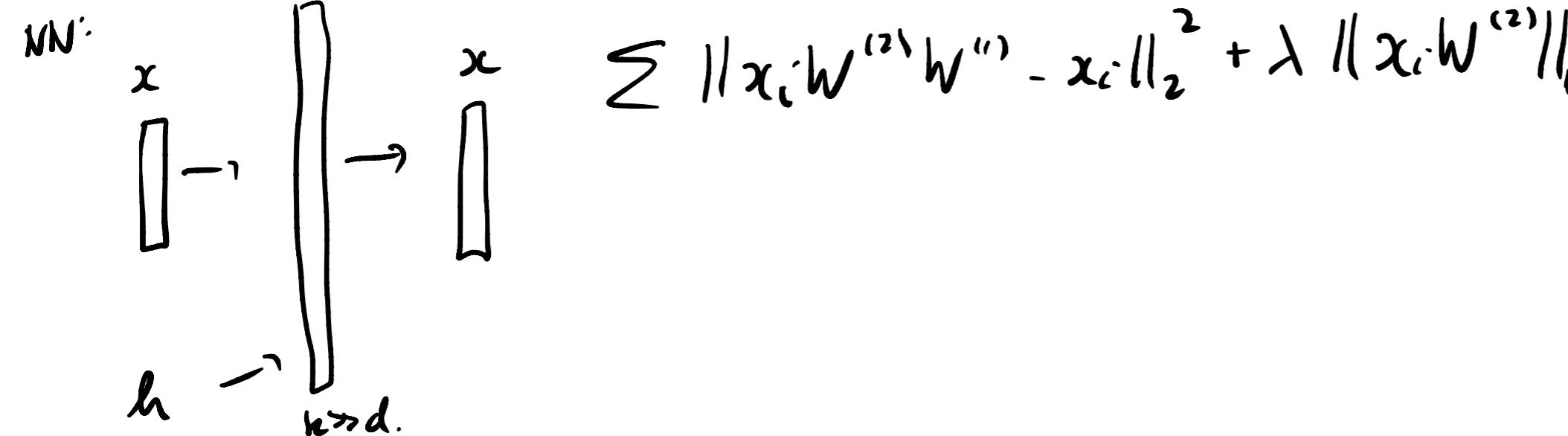
Otherwise typically different, why would we use factorization?

Let's encode specific properties: high-dimensional h that is sparse

h is sparse if a small # of elements are non-zero, $h_i \in \mathbb{R}^k$, $k \gg d$.

$$\min_{D, h_1, \dots, h_n} \sum_{i=1}^n \|x_i - \underbrace{h_i \cdot D}_{\cancel{x_i} \neq x_i}\|_2^2 + \lambda \sum_{i=1}^n \|h_i\|_1 + \lambda_0 \sum_{j=1}^k \|d_j\|_2^2$$

Sparse coding

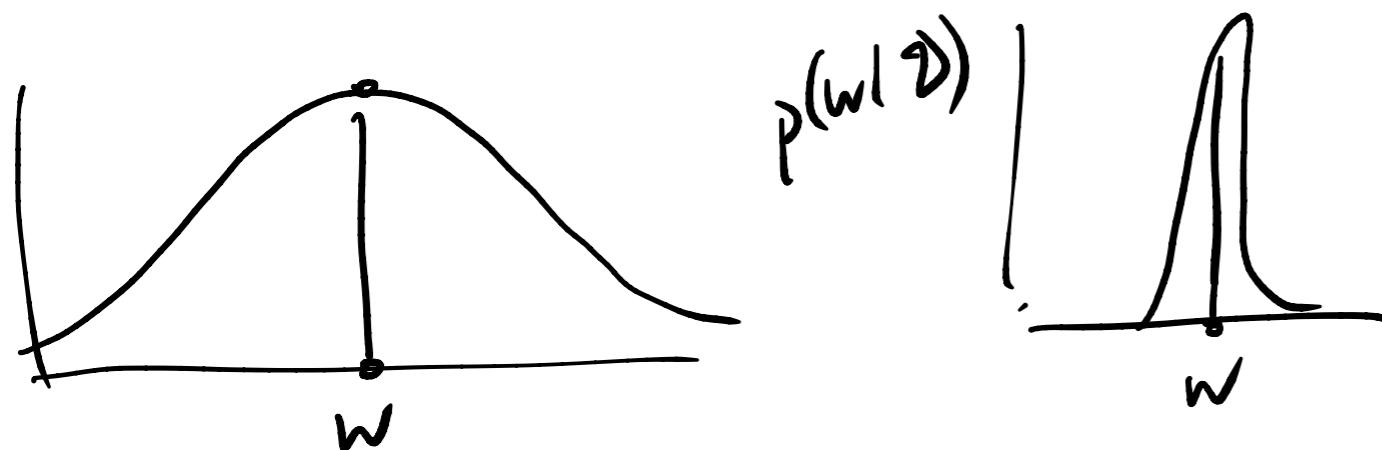


Nov. 28, 2019: Bayesian linear regression

$$p(y|x) = \mathcal{N}(\mu = x^T w, \sigma^2) \quad \mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

$$p(w) = \mathcal{N}(0, \lambda^{-1} I)$$

MAP: $\arg\max_w p(w|\mathcal{D})$ Instead, if we're Bayesian
maintain $p(w|\mathcal{D})$



$$p(w|\mathcal{D}) = \frac{p(\mathcal{D}|w) p(w)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|w) p(w)}{\int p(\mathcal{D}|w') p(w') dw'}$$

$$p(w|\mathcal{D}) = \mathcal{N}(w|m, S) \quad m = (X^T X + \lambda \sigma^2 I)^{-1} X^T y$$

$$S = \sigma^2 (X^T X + \lambda I)^{-1}$$

Conjugate prior: prior for w s.t. that the $p(w|\mathcal{D})$ is of the same form
with different parameters

$$f(x, w) = x^T w$$

$$\bar{f}(x) = E[f(x, w)] \stackrel{\text{random variable}}{=} \int_{\mathbb{R}^d} f(x, w) p(w | \mathcal{D}) dw = m^T x$$

$$\text{Var}(f(x, w)) = \int (f(x, w) - \bar{E}[f(x, w)])^2 p(w | \mathcal{D}) dw \\ = x^T S x$$

If we had picked $p(w)$ = Laplace distr., $p(w | \mathcal{D})$?

- $p(y) \cdot y \in \{0, 1\}$, parameter $w \in [0, 1]$

Bernoulli Bernoulli parameter $\alpha - 1$ successes
 $\beta - 1$ failures

Beta distribution Beta(α, β)

$$\mathcal{D} = \{y_i\}_{i=1}^n, y_i \in \{0, 1\}$$

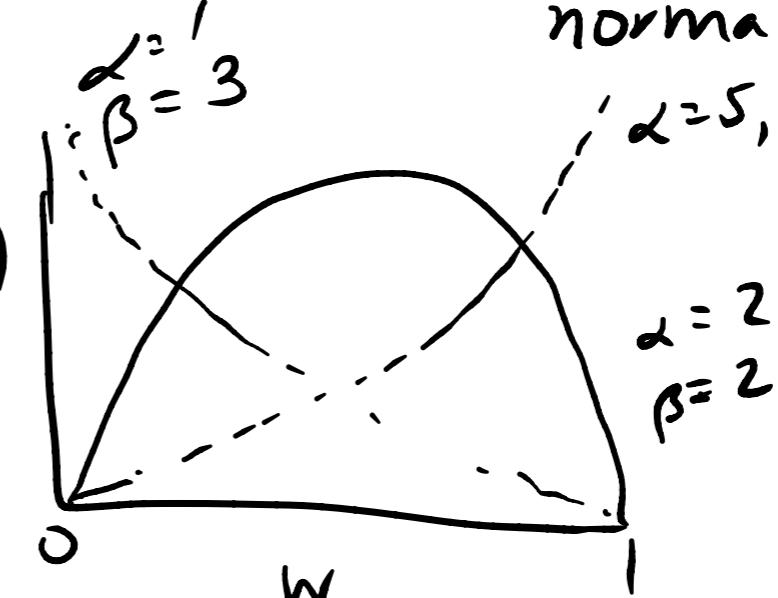
$$p(y) \text{ Bernoulli}(w)$$

$$p(w) \text{ Beta}(\alpha, \beta)$$

$$p(w | \mathcal{D}) = \text{Beta}\left(\alpha + \sum_{i=1}^n y_i, \beta + n - \sum_{i=1}^n y_i\right)$$

successes # failures

$$p(w | \alpha, \beta) = \frac{w^{\alpha-1} (1-w)^{\beta-1}}{\text{normalization term}}$$



$$p(y | \theta) = \int p(y | w) p(w | \theta) dw \quad \text{Posterior predictive distribution}$$

Exercise: What do you do with new data?

$\{y_i\}_{i=n+1}^{n+m}$ e.g. how does $p(w | \theta)$ change
for Bernoulli + Beta?

$$p(w | \theta) = \text{Beta}\left(\alpha + \sum_{i=1}^{n+m} y_i, \beta + n+m - \sum_{i=1}^{n+m} y_i\right)$$
$$\alpha' + \underbrace{\sum_{i=n+1}^{n+m} y_i}_{\text{new data}} , \beta' + m - \sum_{i=n+1}^{n+m} y_i$$

$$p_n(w) = p(w | \theta_n) = \text{Beta}(\alpha', \beta')$$

$$p_n(w | \theta) = \text{Beta}\left(\alpha' + \sum_{i=n+1}^{n+m} y_i, \beta' + m - \sum_{i=n+1}^{n+m} y_i\right)$$