

Support Vector Machine

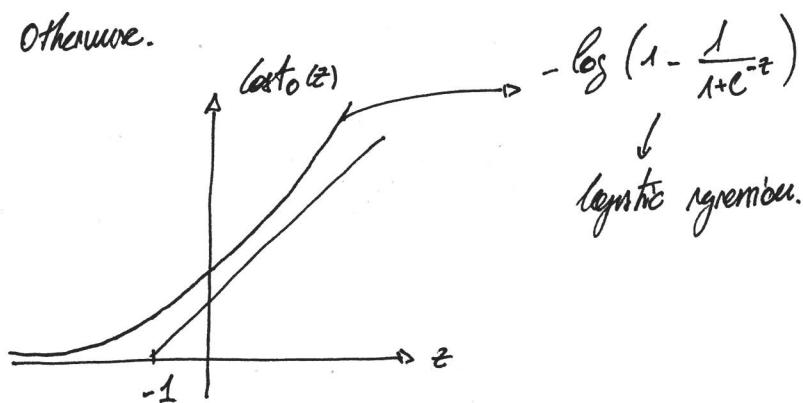
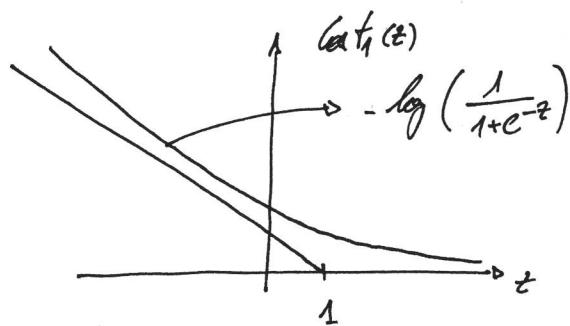
- * Alternative to logistic regression \rightarrow Non-linear solution.

Objective $\rightarrow \min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_2(\theta^T x^{(i)})] \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

\uparrow
gets rid off constant, doesn't impact minimization.

$$\text{cost} = \min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_2(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Hypothesis $\rightarrow h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$

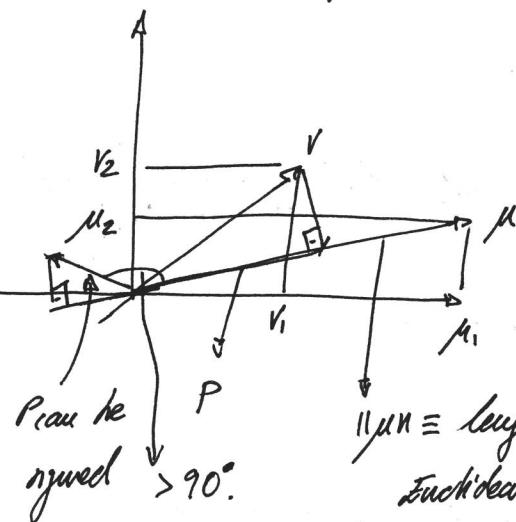


- * to reduce cost \rightarrow If $y=1 \rightarrow \theta^T x \geq 1 \rightarrow \text{cost}_1(z) = 0$.
- * If $y=0 \rightarrow \theta^T x \leq -1 \rightarrow \text{cost}_0(z) = 0$.

- * Theory behind $C = \frac{1}{\lambda}$; trade off between error cost and regularization.

- * If $C \gg 1$, it will cause overfit.

+ Vector inner product \rightarrow



$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}; \nu = \begin{bmatrix} \nu_1 \\ \nu_2 \end{bmatrix}$$

$P \equiv$ length of the projection of vector ν on μ .

$$P \in \mathbb{R}$$

$\|\mu\nu\| \equiv$ length of vector μ $P' \equiv$ length of the projection of vector μ on ν .
Euclidean distance.

$$\|\mu\nu\| = \sqrt{\mu_1^2 + \mu_2^2} \in \mathbb{R}$$

$$\mu^T \cdot \nu = \nu^T \cdot \mu = P \cdot \|\mu\nu\| = P' \cdot \|\nu\| = \nu_1 \mu_1 + \nu_2 \mu_2$$

+ SVM decision boundary \rightarrow

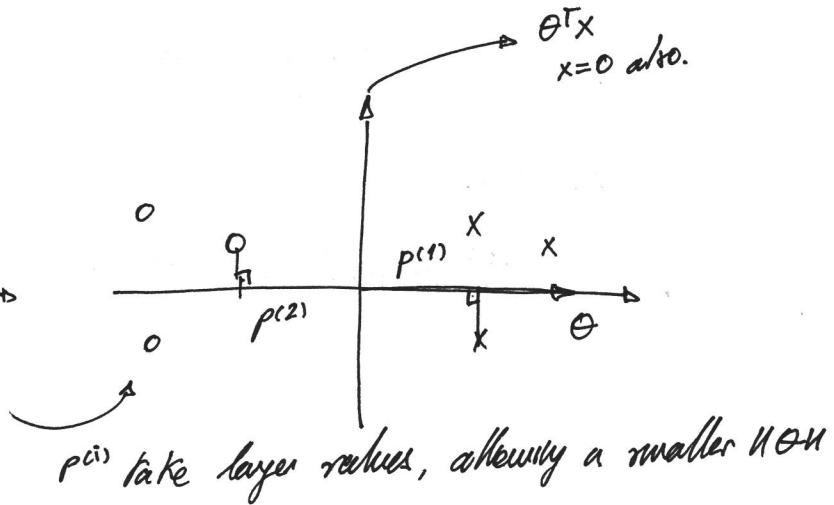
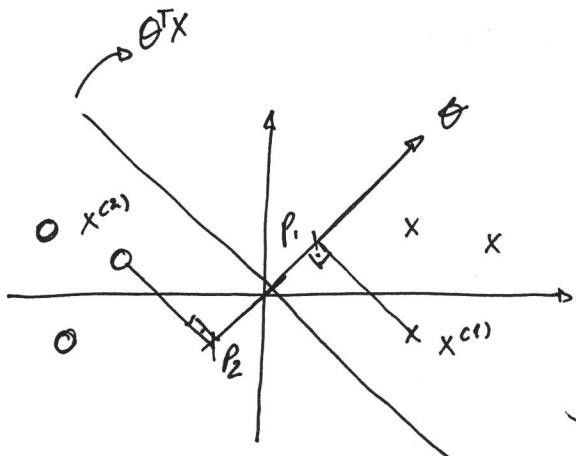
$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n \theta_i^2 \quad \text{if } \theta^T \geq 1 \text{ if } y=1$$

$$\text{if } \theta^T \leq -1 \text{ if } y=0.$$

interpretation $\rightarrow \theta_0 = 0; n=2$

$$\min_{\theta} \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2 \leftarrow \text{minimum lies at origin (0,0) already occupied.}$$

$$\theta^T x^{(i)} = \rho^{(i)} \|\theta\| = \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$

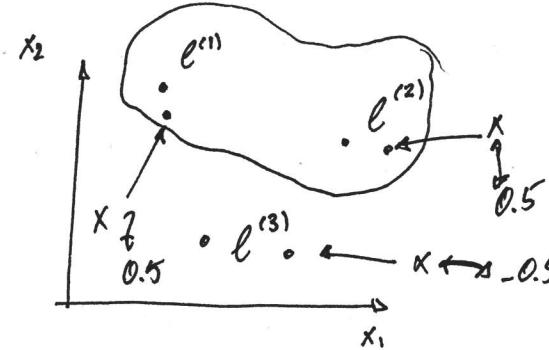


* Kernels →

$$z = \theta^T f = \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \dots + \theta_m f_m$$

$$f_i^{(a)} = \text{similarity}(x^{(a)}, l^{(i)}) ; \quad f \in \mathbb{R}^{m+1} ; \quad \theta \in \mathbb{R}^{m+1} ; \quad f = \begin{bmatrix} f_0^{(a)} \\ f_1^{(a)} \\ \vdots \\ f_m^{(a)} \end{bmatrix}$$

- landmarks are chosen on every input sample → m samples.



Arbitrary 2 dimensions.

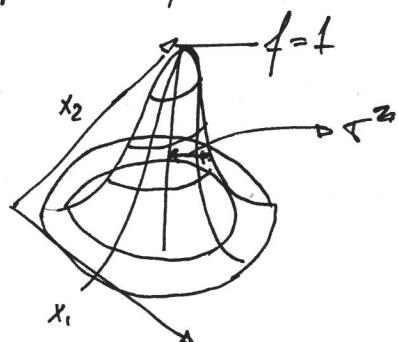
$$\theta_0 = -0.5, \theta_1 = 1; \theta_2 = 1; \theta_3 = 0$$

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0 \rightarrow \hat{y} = 1.$$

Given this θ values, any point close to $l^{(1)}, l^{(2)}$ → $\hat{y} = 1$. Other, $\hat{y} = 0$.

- Gaussian Kernel →

$$f = \text{similarity}(x^{(a)}, l^{(i)}) = \exp\left(-\frac{\|x^{(a)} - l^{(i)}\|^2}{2\tau^2}\right)$$



τ^2 = control how sharp the function is, so for small τ less penalty to the difference.

change dim.

- loss function →

$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1+y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \sum_{j=1}^m \frac{\theta_j^2}{2} \right]$$

* In reality \rightarrow

$$\theta_j^? = -\theta^T \cdot \theta \rightarrow -\theta^T \cdot M \cdot \theta$$

If we have large data set, really expensive to compute. $m=10^6$

* Choosing parameters \rightarrow

$C (= \frac{1}{\alpha}) \rightarrow$ Large $C \equiv$ lower bias, high variance.

Small $C \equiv$ high bias, low variance.

} overfit.

$\tau^2 \rightarrow$ Large $\tau^2 \equiv$ higher bias, low variance

Small $\tau^2 \equiv$ low bias, high variance

- In practice \rightarrow

1. Need to choose C

2. Need to choose Kernel \rightarrow Not all functions work, roughly "Kernel Theorem".
to they don't

A. No Kernel; "Linear Kernel" $\rightarrow y=1 \text{ if } \theta^T x \geq 0.$

Discard

n large and small m \rightarrow avoid overfitting.

B. Gaussian \rightarrow

Need to choose τ^2 .

n small and/or large m.

Important to do feature scaling. \rightarrow Impacts $\|x - \theta\|^2$

* Other Kernels →

- Polynomial Kernel →

$$K(x, l) = (x^T l + \text{const})^{\text{degree}}$$

* Usually need when x, l non negative.

* Need to decide const. degree.

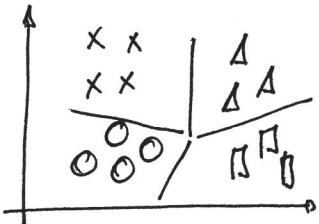
- They, chi-square, histogram intersection.

* Multiclass Classification →

- Use one-vs-all method

- Train K SVM, one to distinguish $y = i$ from the rest $i = 1, 2, \dots, K$

$$\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)} \rightarrow y=1, y=2, \dots, y=K.$$



* Logistic regression vs SVM →

- n is large relative to $m \rightarrow$ logistic regression or SVM linear kernel.

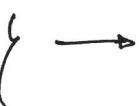
- n is small ($1-1000$), m ($10-10,000$) →

SVM gaussian Kernel

- n is small, m is large →

$$n = 1-1000$$

$$m = 50,000+$$



1. Create new features.

2. Logistic regression or SVM linear kernel.

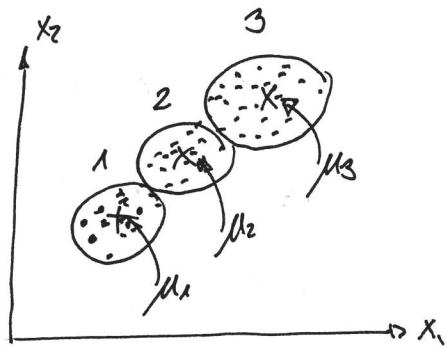
K-Means Clustering

- * Unsupervised learning
- * Algorithm that groups point by itself.
- * Implementation →

Inputs →

$K \equiv$ Number of clusters.

Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$



1. Randomly initialize K cluster centroids $\mu_1, \mu_2, \mu_3, \dots, \mu_K$.

2. Repeat ↴

for $i=1:m$,

Cluster Assignment. → | $c^{(i)} :=$ index (from 1 to K) of cluster centroid to which $x^{(i)}$ is assigned $\rightarrow \min_k \|x^{(i)} - \mu_k\|^2$

Move
centroids. → |

for $k=1:K$,

$\mu_k :=$ average (mean) of points assigned to cluster k .

end;

end; ↴

If there's a centroid with no points → delete.

* Optimization →

$c^{(i)}$ = index of cluster ($1, 2, \dots, K$) to which example $x^{(i)}$ is assigned.

μ_k = cluster centroid K ($\mu_k \in \mathbb{R}^n$)

$\mu_{c^{(i)}}$ = cluster centroid to which $x^{(i)}$ has been assigned.

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

- for each sample minimize the distance between the point and its cluster centroid.

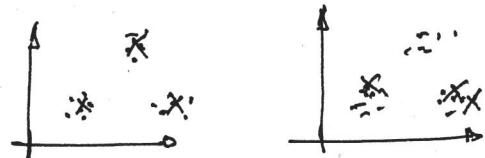
* Random initialization →

- k-Mean may end up in a local optima

- Run several randomizations and pick

clustering lower cost.

- If K is very large, many multiple random initializations won't make much difference.



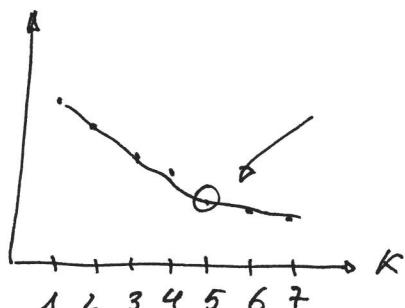
* Number of Clusters →

- Usually chosen manually, though data normalization.

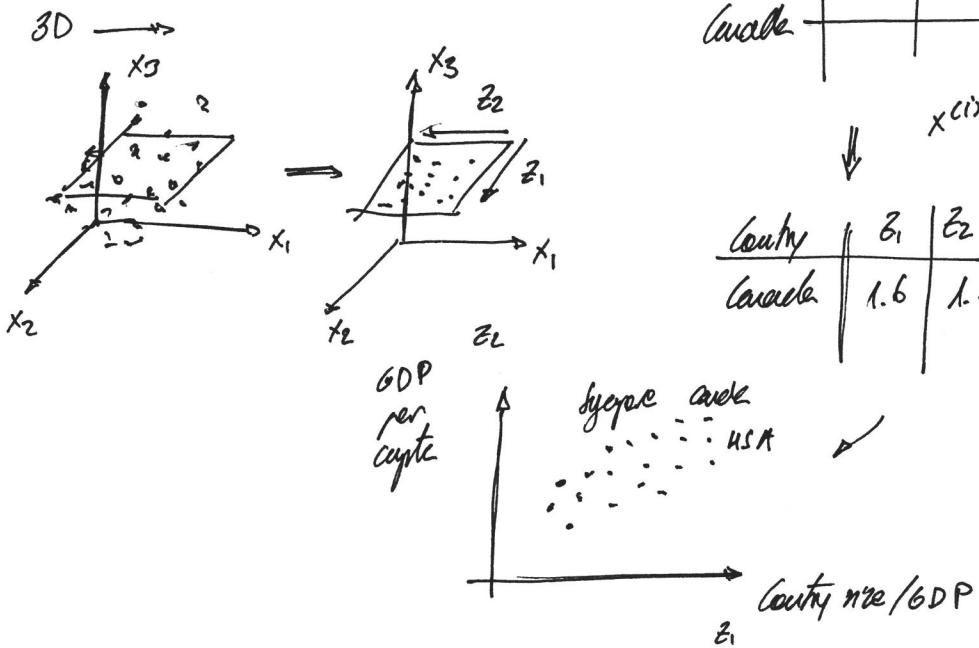
Cost

- Elbow →

- Based on next metric needs.
Elbow is not always clear.



* Data compression & visualization →



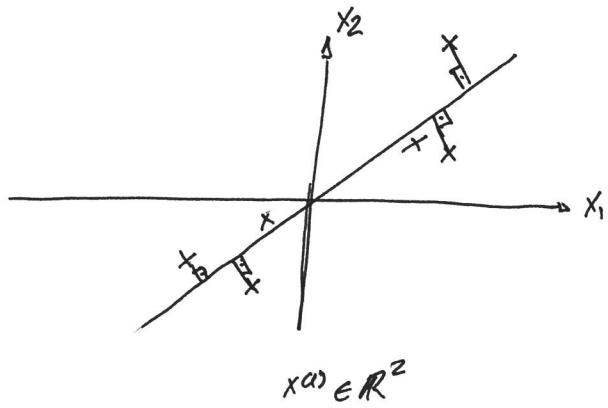
Country	GDP	Per capita GDP	Human Dev index	life expect.	Poverty index
Canada					

$$x^{(1)} \in \mathbb{R}^{50}$$

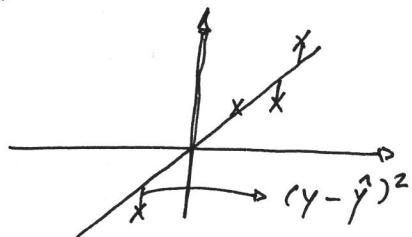
Country	z_1	z_2
Canada	1.6	1.2

$$z^{(1)} \in \mathbb{R}^2$$

Principal Component Analysis



difference to linear regression



- * Finds a lower dimension base in which the projection has the lower distance to the points.

* Algorithm →

A. Data Normalization → Each feature.

$$x_j = \frac{x_j - \mu}{\sqrt{\sigma^2}}$$

B. Reduce data from n-dimension to k-dimension →

$$x^{(i)} \rightarrow z^{(i)} \in \mathbb{R}^k$$

1. Compute "covariance matrix" →

$$\sum = \frac{1}{m} \sum_{i=1}^n x^{(i)} \cdot x^{(i)T} \quad \begin{matrix} \sim \\ \text{nxn} \end{matrix} \quad \Rightarrow \text{ngva} = \frac{1}{m} x^T \cdot x \quad \begin{matrix} \sim \\ \text{nxn} \end{matrix}$$

2. Compute "eigenvectors" of matrix \sum →

"singular value decomposition"

$$[U, S, V] = \text{svd}(\text{ngva}) \quad \begin{matrix} \sim \\ \text{nxn} \end{matrix}$$

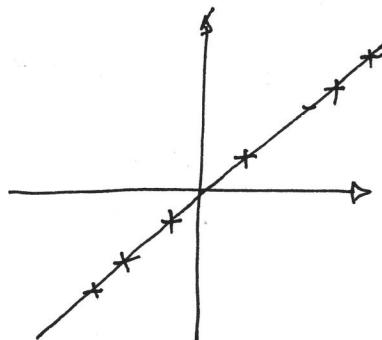
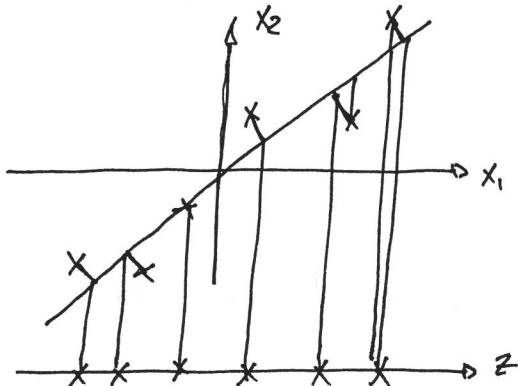
$$U = \begin{bmatrix} | & | & | \\ \mu^{(1)} & \mu^{(2)} & \rightarrow \mu^{(n)} \\ | & | & | \end{bmatrix}$$

K vectors from matrix

$$U_{\text{reduce}} = \dim(n, k) \equiv U(:, 1:k)$$

$$Z = U_{\text{reduce}}^T \cdot X \longrightarrow \dim(m, k)^T \cdot \dim(n \times 1) = \dim(k, 1)$$

* Reconstruction from compressed data \rightarrow



$$X_{\text{approx}} = U_{\text{reduce}} \cdot Z$$

$$Z = U_{\text{reduce}}^T \cdot X$$

* Selecting #PCA \rightarrow

- Objective \rightarrow minimize average square error projection.

$$\text{Error} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$$

- Total variation of my data \rightarrow

On average how far are my samples from 0.

$$\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

* Choose K to be the smallest in which we retain a selected variance \rightarrow

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x^{(i)}_{\text{approx}}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 (1\%) \rightarrow 99\% \text{ of variance retained.}$$

* Algorithm \rightarrow

- Iterate through K values \rightarrow

* Compute U reduce: $x^{(1)} \rightarrow x^{(m)}$

* Calculate $x^{(1)}_{\text{approx}}, \dots, x^{(m)}_{\text{approx}}$.

* Check variance retained.

- K values ($1, 2, \dots, 17$); 17 ready kept it $\rightarrow K=17$.

- How to calculate in Octave \rightarrow

$[U, S, V] = \text{svd}(n \times m)$

$$S = \begin{bmatrix} s_{11} & 0 & \dots & 0 \\ 0 & s_{22} & & \\ \vdots & & \ddots & \\ 0 & \dots & 0 & s_{nn} \end{bmatrix}$$

n \downarrow

For a given $K \rightarrow$

$$1 - \frac{\sum_{i=1}^K s_{ii}}{\sum_{i=1}^n s_{ii}} \leq 0.01$$

* PCA applications \rightarrow

1. Speed-up supervised learning \rightarrow

* Dimensionality reduction.

$$(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)}) \in \mathbb{R}^{10000}$$

PCA

$$(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)}) \in \mathbb{R}^{1000}$$

* Mapping should only be applied by the training, validation or test sets.

2. Data visualization $\rightarrow k=2, k=3$.

3. Bad we \rightarrow Percent encyclicity.

* Doing this dimensionality reduction may mean on losing relationships

between dimensions reduced and labels.

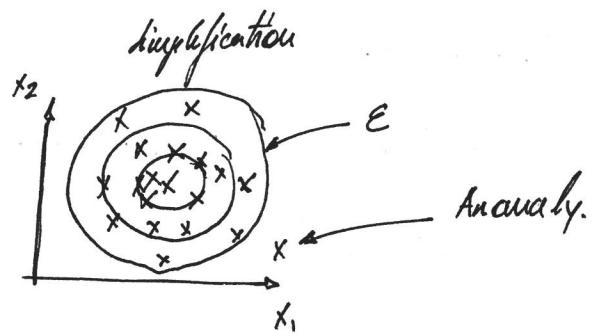
Anomaly Detection

$$\text{Dataset} = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$$

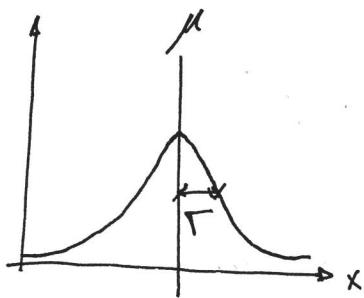
$$\text{Model } p(x) ; x = (x_1, x_2, \dots, x_n)$$

$$\rightarrow p(x_{\text{test}}) < \epsilon \rightarrow \text{Anomaly.}$$

$$p(x_{\text{test}}) \geq \epsilon \rightarrow \text{OK.}$$

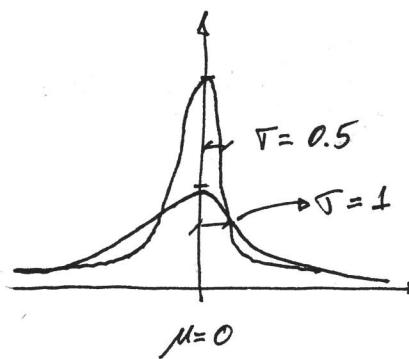


* Gaussian distribution \rightarrow



Probability distribution.

$$x \sim N(\mu, \sigma^2)$$



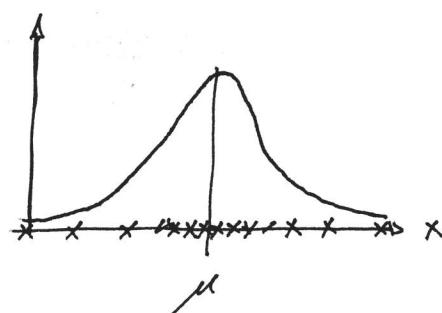
$$p(x) = p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- Parameter estimation \rightarrow

$$\text{Dataset} = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} ; x^{(i)} \in \mathbb{R}.$$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$



maximum likelihood estimates.

* Algorithm \rightarrow

Training set $\rightarrow \{x^{(1)}, \dots, x^{(m)}\}$

$x^{(i)} \in \mathbb{R}^n$

$$p(x) = p(x_1; \mu_1, \sigma_1^2) \cdot p(x_2; \mu_2, \sigma_2^2) \cdot \dots \cdot p(x_n; \mu_n, \sigma_n^2)$$

Independent assumption.

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

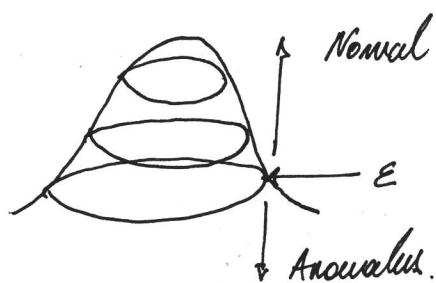
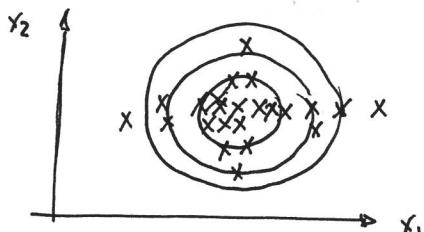
$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} ; \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

1. Choose features that might indicate anomalous examples.

2. Fit parameters $\rightarrow \mu_i, \sigma_i^2$

3. Create model $\rightarrow p(x)$

4. Anomaly if $p(x) < \epsilon$.



- Algorithm evaluation →

* the labeled data for evaluation. Eg. →

- Training set → 6000 good → $p(x)$ model.

- Cross-validation → 2000 good / 10 anomalies.

- Test → 2000 good / 10 anomalies.

* Flow →

1. Fit model $p(x)$ on training set → $\{x^{(1)}, \dots, x^{(m)}\}$

2. On cross-validation / test example x →

$$y = 1 \text{ if } p(x) \leq \epsilon \text{ Anomaly}$$

$$y = 0 \text{ if } p(x) > \epsilon \text{ Normal.}$$

3. Possible evaluation metrics → Very skewed data.

* True positive, false positive, false negative, false positive.

* Precision / Recall.

* F_1 -Score → harmonic mean of precision / recall.

4. Optional → You can use cross-validation to set parameter ϵ .

- When to use instead of supervised learning →

* Very small number of positive examples, and large negative.

* Many different types of anomalies

* Future anomalies may look nothing like previous.

- Choosing features \rightarrow

* Features are modeled as gaussians

* If features don't have this distribution, transformations \rightarrow

$$* \log(x + c)$$

$$* x^c$$

* Anomaly not captured by current features \rightarrow Create new feature.

Initial features $\rightarrow x_1 \equiv$ memory use of computer.

$x_2 \equiv$ number of disk accesses/sec

$x_3 \equiv$ CPU load.

$x_4 \equiv$ Network traffic

New feature $\rightarrow x_5 = \frac{\text{CPU load}}{\text{Network traffic}}$ or $\frac{\text{CPU load}^2}{\text{Network traffic}}$.

This will capture anomalies of computer stuck in ref. loops.

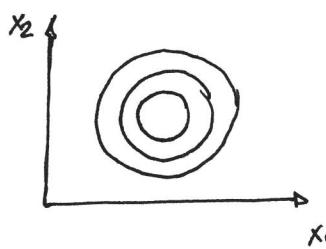
- Multivariate Gaussian Distribution \rightarrow

$$f(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} ((x-\mu)^T \Sigma^{-1} (x-\mu))}$$

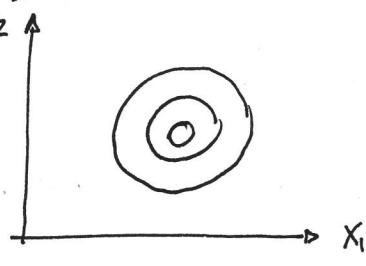
$$\mu \in \mathbb{R}^n$$

$\Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix). ; $|\Sigma| \equiv$ determinant of covariance.

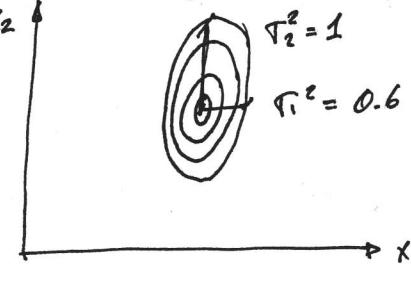
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



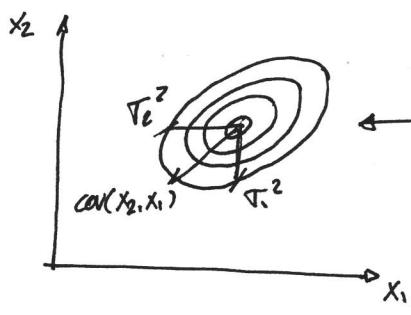
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$



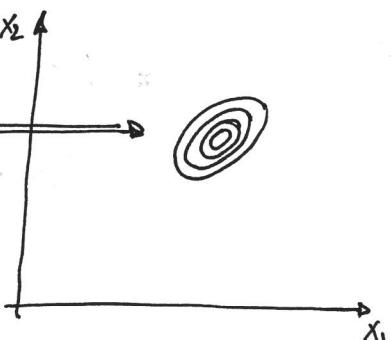
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$



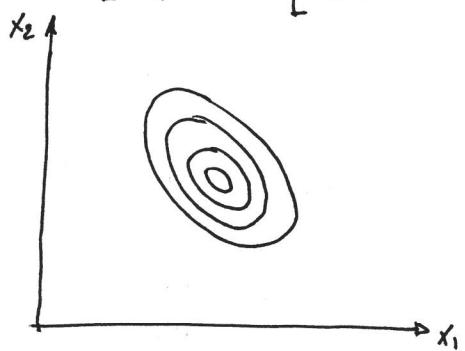
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$



- * Differently from covariance with itself. covariance between ~~values~~ peaks while inaccuracy.

- * Application →

$$\text{Training set} = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad ; \quad \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

- Given example → $p(x_{\text{test}}) < \epsilon?$

- * Physical model is a special case of the multivariate model →

It assumes that variables are independent.

$$\Sigma = \text{diag matrix} = \begin{bmatrix} \tau^2 & 0 & \cdots & 0 \\ 0 & \tau_2^2 & \cdots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \tau_n^2 \end{bmatrix}$$

Physical Model

- * Manually create features to highlight anomalies

$$x_3 = \frac{x_1}{x_2}$$

- * Computationally cheaper

- * Work ok with small number of samples m .

- Features that are not linearly independent $\rightarrow x_3 = x_1 + x_5$

Σ^{-1} not invertible.

Mathematical Model

- * Automatically captures correlation between features.

* More expensive \rightarrow

$$\Sigma \in \mathbb{R}^{n \times n}; \Sigma^{-1}$$

* Must have $m > n$ or else

Σ is not invertible

Σ nuclear matrix.

Recommender Systems \rightarrow

* Latent Bond Recommendation \rightarrow

Some Features.

Movies	Person 1	Person 2	Person 3	Person 4	x_1 (Revenue)	x_2 (action)
Movie 1	5	5	0	0	0.9	0
Movie 2	5	?	?	0	1.0	0.01
Movie 3	?	4	0	?	0.99	0
Movie 4	0	0	5	4	0.1	1.0
Movie 5	0	0	5	?	0	0.9

- Notation \rightarrow

* $n_u \equiv \# \text{ users.}$

$$\rightarrow n_u = 4; n_m = 5.$$

* $n_m \equiv \# \text{ movies.}$

* $r_{ui,j} = 1 \text{ if user } j \text{ has reviewed movie } i.$

* $y_{ci,j} = \text{Rating given by user } j \text{ to movie } i.$

defined only if $r_{ui,j} = 1.$

* $\theta^{(j)}$ = Parameter vector for user j .

* $x^{(i)}$ = Parameter vector for movie i . Features.

* Prediction $\rightarrow (\theta^{(j)})^T (x^{(i)})$

* $m^{(j)} \equiv \# \text{ of movies rated by user } j.$

- Cost function \rightarrow

$$J(\theta^{(1)}, \dots, \theta^{(n)}) = \min_{\theta^{(1)}, \dots, \theta^{(n)}} \frac{1}{2} \sum_{j=1}^m \sum_{i: r(i,j) \neq 1} (\theta^{(j)\top} x^{(i)} - y^{(i,j)})^2$$

$$+ \frac{\lambda}{2} \sum_{j=1}^m \sum_{k=1}^n (\theta_k^{(j)})^2$$

Remove $m(j)$ because a constant factor doesn't impact the optimization value for $\theta^{(j)}$

- Gradient descent update \rightarrow

$$\theta_k^{(j)} \leftarrow \theta_k^{(j)} - \alpha \sum_{i: r(i,j)=1} (\theta^{(j)\top} x^{(i)} - y^{(i,j)}) \cdot x_k^{(i)} \quad \text{for } k=0.$$

$$\theta_k^{(j)} \leftarrow \theta_k^{(j)} - \alpha \left[\sum_{i: r(i,j)=1} (\theta^{(j)\top} x^{(i)} - y^{(i,j)}) \cdot x_k^{(i)} + \lambda \theta_k^{(j)} \right] \quad \text{for } k \neq 0.$$

* Collaborative Filtering \rightarrow

- In this case, no knowledge of the features \rightarrow

Movie	Person 1	Person 2	Person 3	Person 4	x_1 (relevance)	x_2 (action)
Movie 1	5	5	0	0	?	?
Movie 2	5	?	?	0	?	?
Movie 3	2	4	0	?	?	?
Movie 4	0	0	5	4	?	?

- If we know $\theta^{(j)}$ \rightarrow

$$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \rightarrow (\theta^{(1)})^\top x^{(1)} = 5 \rightarrow \text{We can infer features.}$$

- cost function \rightarrow

$$J(x^{(1)}, \dots, x^{(n)}) = \min_{\theta^{(1)}, \dots, \theta^{(n)}} \frac{1}{2} \sum_{j: r(i,j)=1} (\theta^{(j)T} x^{(i)} - y^{(i),j})^2 + \frac{1}{2} \sum_{k=1}^m \sum_{j=1}^n (x_k^{(i)})^2$$

- Gradient descent update \rightarrow

$$x_k^{(i)} \leftarrow x_k^{(i)} - \alpha \left(\sum_{j: r(i,j)=1} (\theta^{(j)T} x^{(i)} - y^{(i),j}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

* Recommender Algorithm \rightarrow

- cost function \rightarrow

$$J(x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}) = \frac{1}{2} \sum_{c(i,j): r(i,j)=1} (\theta^{(j)T} x^{(i)} - y^{(i),j})^2 + \frac{1}{2} \sum_{i=1}^m \sum_{k=1}^n (x_k^{(i)})^2 + \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \theta_k^{(j)}^2$$

* $c: r(i,j)=1 \rightarrow$ For every movie rated by user

$$x_0 = \theta_0 = 0$$

* $j: r(i,j)=1 \rightarrow$ For every user that rated movie i .

$$x \in \mathbb{R}^n, \theta \in \mathbb{R}^n$$

- Algorithm \rightarrow

1. Initialize $x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}$ random variables.

2. Minimize $J(x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)})$

$$x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j: r(i,j)=1} (\theta^{(j)T} x^{(i)} - y^{(i),j}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

$$\theta_k^{(i)} := \theta_k^{(i)} - \alpha \left(\sum_{c(i,j)=1} (\theta^{(j)T} x^{(i)} - y^{(i),j}) x_k^{(i)} + \lambda \theta_k^{(i)} \right)$$

3. For user with parameters θ and movie with learned features x , predict a rating $\theta^T x$.

- Low Rank Matrix Factorization \rightarrow

* Collaborative filtering \rightarrow

Movie	Person 1	Person 2	Person 3	Person 4
Movie 1	5	5	0	0
Movie 2	5	?	?	0
Movie 3	?	4	0	?
Movie 4	0	0	5	4
Movie 5	0	0	5	?



$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \quad n_{mu} \equiv \# \text{ movies.}$$

$n_{ui} \equiv \# \text{ users.}$

$$y \in \mathbb{R}^{(n_{mu}, n_u)}$$

$$Y = \begin{bmatrix} \theta^{(1)T} x^{(1)} & \theta^{(2)T} x^{(2)} & \dots & \theta^{(n_u)T} x^{(1)} \\ \vdots & & & \\ \theta^{(1)T} x^{(n_m)} & \theta^{(2)T} x^{(n_m)} & \dots & \theta^{(n_u)T} x^{(n_m)} \end{bmatrix} \quad \in \mathbb{R}^{(n_{mu}, n_u)}$$

$$x = \begin{bmatrix} \cdots x^{(1)T} \cdots \\ \vdots \\ \cdots x^{(n_m)T} \cdots \end{bmatrix} ; \quad \Theta = \begin{bmatrix} \cdots \theta^{(1)T} \cdots \\ \vdots \\ \cdots \theta^{(n_u)T} \cdots \end{bmatrix} ; \quad x \in \mathbb{R}^{(n_m, \text{ movie features})}$$

$; \quad \Theta \in \mathbb{R}^{(n_u, \text{ movie features})}$

$$y = x \Theta^T$$

- Find related movies \rightarrow How to find movie j related to movie i ?

$$\|x^{(i)} - x^{(j)}\| \rightarrow \text{small.}$$

* mean normalization \rightarrow

- then without rated movies.

- Problem \rightarrow

* If we run this over the cost function $\theta^{(new user)} = 0$.

* hence regularization for $\theta^{(j)}$ will keep $\theta^{(j)}$ value low,
and there's no defined targets $\rightarrow \theta^{(j)} = 0$.

$$y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix} \rightarrow \mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 3 \\ 2.25 \\ 2.25 \end{bmatrix} \rightarrow y - \mu$$

- For user j , on movie i predict \rightarrow

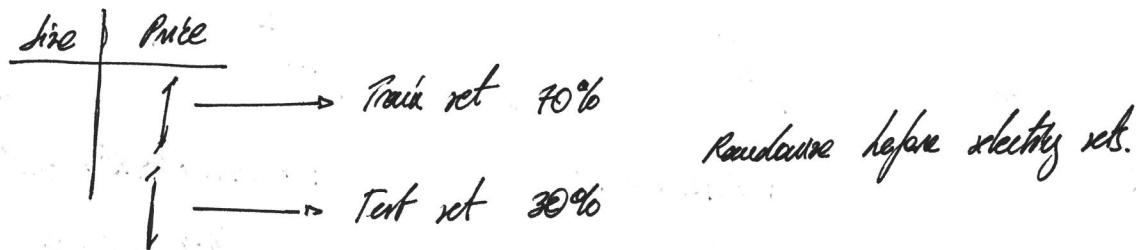
$$\theta^{(j)T} x^{(i)} + \mu_i$$

* Haer 5 (Ever) \rightarrow New user!

$$\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rightarrow \underbrace{(\theta^{(5)T} x^{(i)})}_{L=0} + \mu_i //$$

Applying machine learning

- * Evaluating hypothesis \rightarrow



- On classification you may use other error definitions instead of cross-entropy \rightarrow

Classification
error

$$\text{error}(h_\theta(x), y) = \begin{cases} 1 & \text{if } h_\theta(x) \geq 0.5, y=0 \text{ or} \\ & h_\theta(x) < 0.5, y=1 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Test error} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \text{error}(h_\theta(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

- * Model classification \rightarrow

Train set 60%

$$d=1 \rightarrow h_\theta(x) = \theta_0 + \theta_1 x \rightarrow \theta^{(1)}$$

Cross-validation 20%

$$d=2 \rightarrow h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

Test set 20%

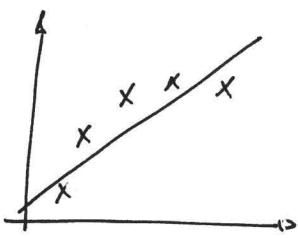
$$\vdots$$

$$d=10 \rightarrow h_\theta(x) = \theta_0 + \theta_1 x^2 + \dots + \theta_{10} x^{10}$$

- The cross-validation set to select the model \rightarrow

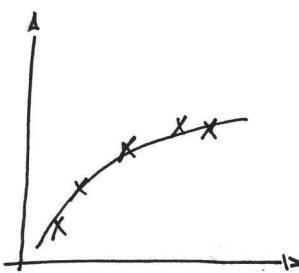
$$J(\theta^{(1)}), \dots, J(\theta^{(10)}) \rightarrow \text{we called error}$$

* Diagnose bias vs variance \rightarrow Model selection.



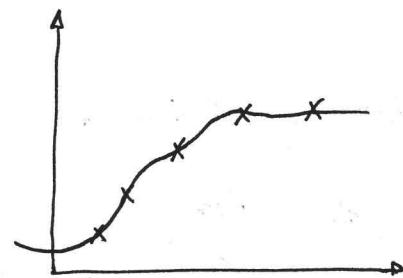
$$\theta_0 + \theta_1 x$$

High bias
(underfit)



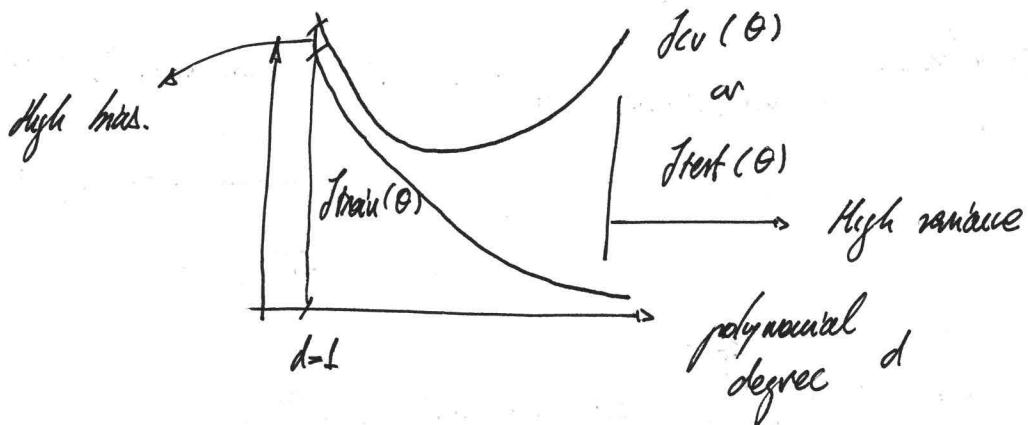
$$\theta_0 + \theta_1 x + \theta_2 x^2$$

Correct.



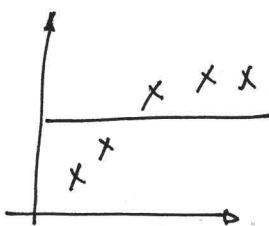
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

High variance
(overfit).



* Regularization and bias/variance \rightarrow

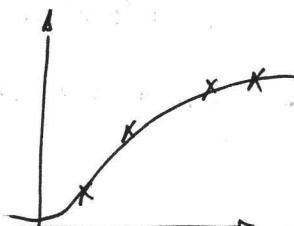
$$\text{Model} \rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4; \text{Training cost} \rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$



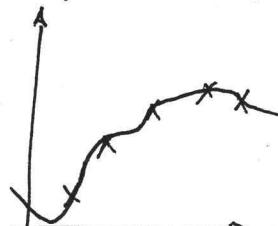
Large λ
High bias

$$\lambda = 10000$$

$$\theta_0 \approx \theta_1 \approx \theta_2 \approx \dots \approx 0$$

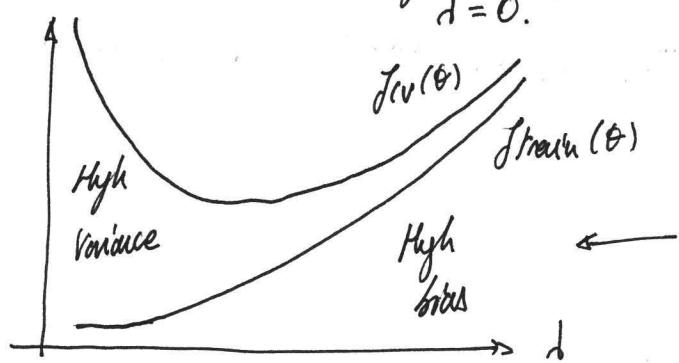


Intermediate λ
Correct



Small λ
High Variance
 $\lambda = 0$.

$$J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2; J_{\text{cv}}(\theta); J_{\text{test}}(\theta)$$

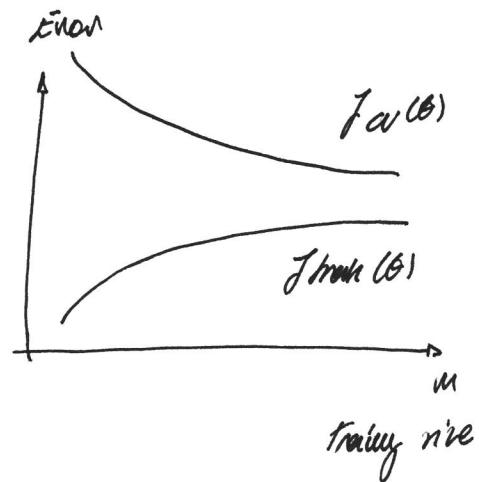


Select one with
lowest cost.

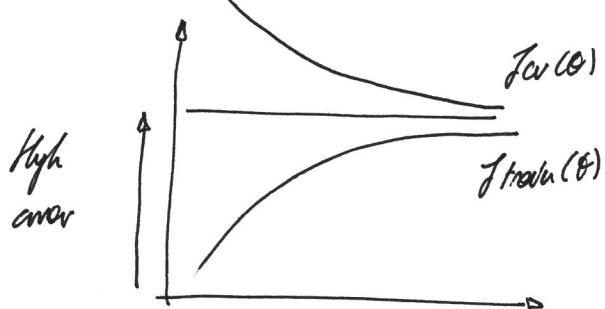
* Learn my aunc →

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J_{\text{cv}}(\theta) = \frac{1}{m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

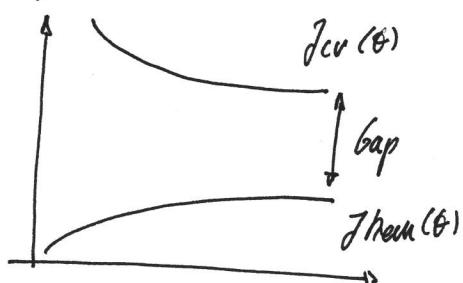


- High bias →



- + Simple model
- * More data won't help.

- High variance →



- * More data will likely help.

* What to do next →

1. More training examples → High variance.

2. Smaller set of features → High variance.

3. Additional features → High bias.

4. Polynomial features → High bias

5. Decreasing d → High bias.

6. Increasing d → High variance.

