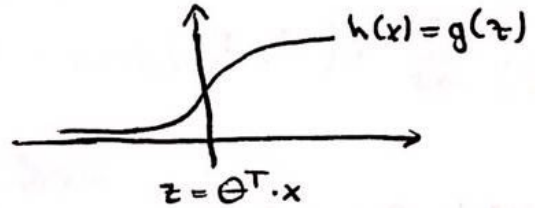


There is one more powerful and very widely used learning algorithm: Support Vector Machine (SVM): compared to both logistic regression and neural networks, SVM sometimes gives a cleaner and more powerful way of learning complex non-linear functions.

Logistic Regression ile başlayıp değışiklikler ile SVM'e anlagayacağız.

$$h(x) = \frac{1}{1 + e^{-\Theta^T \cdot x}}$$



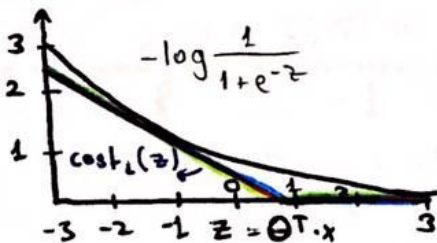
- If  $y=1$ , we want  $h(x) \approx 1$ ,  $\Theta^T \cdot x \gg 0$  ki hipotez doğru olsun.
- If  $y=0$ , we want  $h(x) \approx 0$ ,  $\Theta^T \cdot x \ll 0$

Cost function tüm training ex.s ile hesaplanıyordu. Tek bir training ex.'in etkisine (cost'una) bakalım:

$$- (y \log h(x) + (1-y) \log (1-h(x)))$$

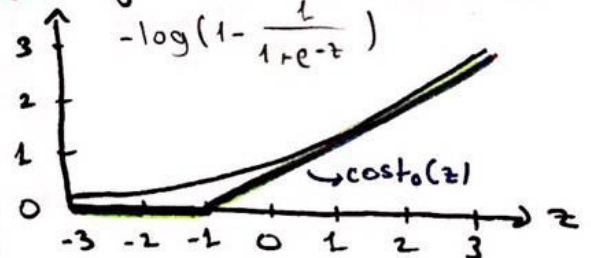
$$- y \log \frac{1}{1+e^{-\Theta^T \cdot x}} - (1-y) \log \left(1 - \frac{1}{1+e^{-\Theta^T \cdot x}}\right)$$

If  $y=1$  (want  $\Theta^T \cdot x \gg 0$ )



$y=1$  ise  $z$  büyüdükçe ilgili training ex. cost'u düşüyor üstteki  $h(x)$  grafiği de aynı şeyi söyler.  $z$  büyüdükçe  $h(x) \approx 1$  olur. Cost düşer.

If  $y=0$  (want  $\Theta^T \cdot x \ll 0$ ):



**★ SVM için başka bir cost function kullanacağız.** Yukarıda ~~maximize~~ ile çözüldü. It is a pretty close approximation to the cost function used by logistic regression. Slope çok önemli değil. Bu yeni cost function SVM'e computational advantages verir. Give us an easier optimization problem.

- Cost function for logistic regression:

$$\min_{\theta} \underbrace{\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \cdot \underbrace{(-\log h(x^{(i)}))}_{\text{SVM için burası, cost}_1(\theta^T x^{(i)})} + (1-y^{(i)}) \cdot \underbrace{(-\log(1-h(x^{(i)})))}_{\text{Burası da cost}_0(\theta^T x^{(i)})} \right]}_{\text{Term A}} + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2}_{\text{Term B}}$$

- Cost Function for Support Vector Machine:

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m y^{(i)} \cdot \text{cost}_1(\theta^T \cdot x^{(i)}) + (1-y^{(i)}) \cdot \text{cost}_0(\theta^T \cdot x^{(i)}) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

By convention we write things slightly different

- get rid of the  $1/m$  terms. Somewhat but constant  $\min \theta$  doesn't matter.

- $\lambda + 2B$ : Burada  $\lambda$  ile underfitting ( $\lambda \uparrow$ ) vs overfitting ( $\lambda \downarrow$ ) trade-off'u kontrol ediliyordu. Şimdi convention olarak  $\lambda$  yerine  $C$  kullanılarak ve:  $C \lambda + B$  kontrol edilecek.  $\lambda$  ile ters etki yapacak ama mantık aynı!  $C = \frac{1}{\lambda}$  gibi düşünebiliriz. ise aynı  $\min \theta$  bulunur!

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

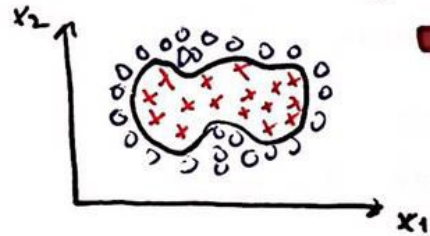
- Unlike logistic regression, the SVM doesn't output a probability. It just makes a prediction of  $y$  being equal to one or zero, directly.

$$h(x) = \begin{cases} 1 & \text{if } \theta^T \cdot x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



• We will start adapting SVMs in order to develop complex nonlinear classifiers. The main technique for doing that is something called kernels.

• If we have a training set as below and we have to find a nonlinear decision boundary to distinguish the positive and negative examples.



• One way to do so is complex polynomial features.

- Predict  $y=1$  if  $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0$
- Predict  $y=0$  otherwise

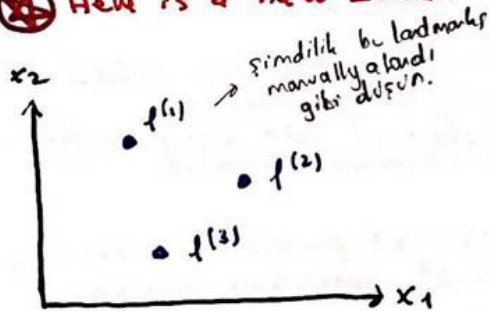
• Yukarıdaki  $\theta^T \cdot x$  i şu şekilde de yazabiliriz:

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots \quad \text{where } f_1 = x_1, f_2 = x_2, f_3 = x_1 \cdot x_2, \dots$$

★ The question is: Is there a better choice of features than high order polynomials? Because it is not clear that these high order polynomials are what we want, and we have seen that when  $n$  (feature #) is large (for ex for a computer vision problem) using high order polynomials become very computationally expensive. There will be too many features.

• So is there a better choice of these features  $f_1, f_2, f_3, \dots$ ?

★ Here is a new Idea: **KERNEL**. Let's define only 3 new features:



• Given the training example  $x$ , compute the new feature depending on proximity to the related landmark  $l^{(i)}$ .  $\|x - l^{(i)}\| \rightarrow$  Euclidean distance b/w point  $x$  and  $l^{(i)}$

$$\rightarrow \text{Given } x: f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$\rightarrow \text{Given } x: f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

$$\rightarrow \text{Given } x: f_3 = \text{similarity}(x, l^{(3)}) = \exp(\dots)$$

I am ignoring  $\sigma$  for the purpose of these slides

★ The similarity function is called Kernel Function and the specific kernel I'm using here is a Gaussian Kernel. Başka bir kernel de kullanılabildi. Bu kullandığın Gaussian kernel.  $k(x, l^{(i)})$

Her training ex için farklı bir feature vektörü olur.  $(f_1, f_2, f_3, \dots)$



### Kernels and Similarity:

$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{2\sigma^2}\right)$

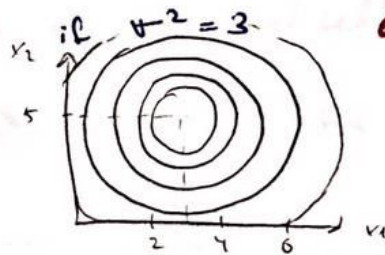
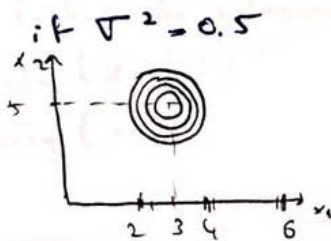
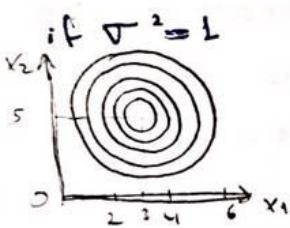
• If  $x \approx l^{(1)}$ :  $f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$

• If  $x$  is far from  $l^{(1)}$ :  $f_1 \approx \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$

- Each landmark defines a new feature:  $l^{(1)} \rightarrow f_1, l^{(2)} \rightarrow f_2, l^{(3)} \rightarrow f_3, \dots$
- Bır bir training ex verince landmarks sayısı kadar feature b-tilir.

Let's plot and try to understand this exponential Gaussian Kernel.

Ex:  $l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$

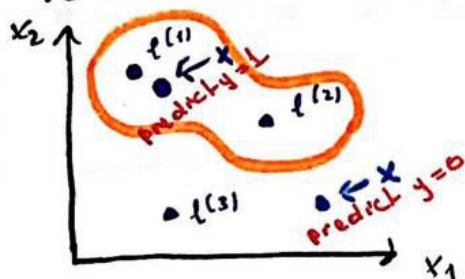


•  $\sigma$  küçüldükçe hızlı sıvırlaşan dık bir dağ, aksi halde düşük eğimli geniş bir dağ grıbr düşün.

Let's see what sorts of hypotheses we can learn:

→ Let's say we have trained our algorithm and our hypothesis: **predict 1** when  **$\Theta_0 + \Theta_1 f_1 + \Theta_2 f_2 + \Theta_3 f_3 \geq 0$**  where  **$\Theta_0 = -0.5, \Theta_1 = 1, \Theta_2 = 1, \Theta_3 = 0$**

• Verilen training ex. için features  $f_1, f_2, f_3$  hesaplanıp prediction yapılacaktır ve kabaca decision boundary çizilecek (intuition için)



• Non training ex,  $x$  baz alınarak:

- $f_1 \approx 1$  since  $x$  is close to  $l^{(1)}$
- $f_2 \approx 0, f_3 \approx 0$  since  $x$  is far away from  $l^{(2)}$  and  $l^{(3)}$
- $\Theta^T \cdot f \Rightarrow \Theta_0 + \Theta_1 = 0.5 \geq 0$  Predict  $y = 1$

• Mavi training ex,  $x$  baz alınarak:

- $f_1 \approx 0, f_2 \approx 0, f_3 \approx 0$
- $\Theta^T \cdot f \Rightarrow \Theta_0 = -0.5 < 0$  thus predict  $y = 0$

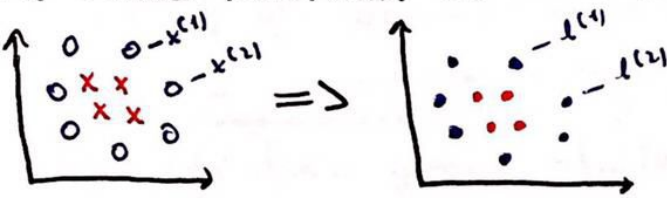
• Bunu devam ettirmek için  $l^{(1)}$  ve  $l^{(2)}$  ya yakın training ex. için  $y = 1$  tahmini yapabiliriz, uzak olanlar için  $y = 0$  tahmini yapılır.

• Bu şekilde bir non-linear decision boundary oluşturabiliriz.

- Kernels ile ilgili bazı detaylardan bahsedeceğiz. Prestiğe nasıl kullanılabileceğine değineceğiz. SVM den bias-variance trade off'ına değineceğiz.

### How to place landmarks?

- For every training example we are going to put a landmark as exactly the same locations as the training examples.



- I'm gonna end up with  $m$  landmarks  $l^{(1)}, l^{(2)}, \dots, l^{(m)}$

- Renkler anlatsın diye var aslında örensiz.

### SVM with Kernels

- Given  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$
- choose  $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$

- Given example  $x$  (it can be in training set, cv set or test set):

$$\begin{aligned} f_1 &= \text{similarity}(x, l^{(1)}) \\ f_2 &= \text{similarity}(x, l^{(2)}) \\ &\vdots \end{aligned} \Rightarrow f = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_m \end{bmatrix} \text{ where } f_0 = 1$$

### For training example $(x^{(i)}, y^{(i)})$ :

$$\begin{aligned} x^{(i)} \rightarrow f_1^{(i)} &= \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} &= \text{sim}(x^{(i)}, l^{(2)}) \\ &\vdots \\ f_m^{(i)} &= \text{sim}(x^{(i)}, l^{(m)}) \end{aligned} \quad \text{where } f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = 1$$

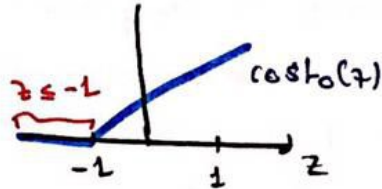
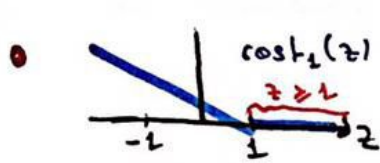
Çünkü training setten gelilen bir  $x^{(i)}$  için kesinlikle bir  $l^{(i)} = x^{(i)}$  olacak bu noktada da  $f_i^{(i)} = 1$  olacaktır.

- We can use  $f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$  instead as  $x^{(i)} \in \mathbb{R}^{n+1}$



• Sometimes people talk about SVMs as Large Margin Classifiers.  
Bu derste bunden bahsedeceğimiz ve SVM hypothesis nasıl görününcü konuşacağız.

•  $\min_{\theta} C \sum_{i=1}^m [y^{(i)} \cdot \text{cost}_1(\theta^T \cdot x^{(i)}) + (1-y^{(i)}) \cdot \text{cost}_0(\theta^T \cdot x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$



- ② If  $y=1$ , we want  $\theta^T \cdot x \geq 1$  (Not just  $\geq 0$ )  
 $\text{cost}_1(z)$  is 0 only when  $z \geq 1$
- ③ If  $y=0$ , we want  $\theta^T \cdot x \leq -1$  (Not just  $\leq 0$ )  
 $\text{cost}_0(z)$  is 0 only when  $z \leq -1$

→  $\theta^T \cdot x \geq 0$  veya  $< 0$  olması yetmiyor. Bunu bir safety factor gibi düşünebiliriz. Bu olay extra safety margin factor bulma eden

→ Bu durumun sonuçlarına (SVM için) bakalım:

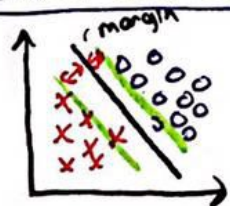
④ Consider a case that  $C$  is very large: Let's say 10 000 (means  $\lambda$  is small. Overfitting)

•  $C$  çok büyükse cost function'ın minimize edilmesi için A termi  $\approx 0$  olmalı.

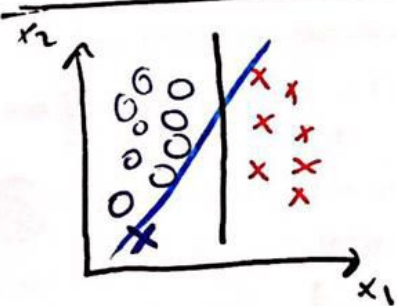
- whenever  $y^{(i)} = 1$  :  $\text{cost}_1(z)$ 'nin 0 olması için  $\theta^T \cdot x^{(i)} \geq 1$  olmalı!
- whenever  $y^{(i)} = 0$  :  $\text{cost}_0(z)$ 'nin 0 olması için  $\theta^T \cdot x^{(i)} \leq -1$  olmalı!

⑤ Sonuçta eğer Term A 0'a yakın ise bu şu demek her  $y^{(i)} = 1$  için  $\theta^T \cdot x^{(i)} \geq 1$  ve her  $y^{(i)} = 0$  için  $\theta^T \cdot x^{(i)} \leq -1$  sağlanmış demek.  
(Term A = 0 ise hepsi sağlanmış demek yakın 0'a kadar kaçan olabilir)  
Bu durum ise decision boundary'e şu şekilde yansır:

SVM Decision Boundary: Linearly Separable Case



• Siyah decision boundary high margin'e sahip. Yani ayırdığı kumeler ıstabilliği fazla. iyi ama bu zaten logistic regression için de eğer cost function = 0 ise yaşanmaz mı? Yine large margin olur



•  $C$  is very large kabul ediyoruz. Vardaki anekdot datalar siyah çizgi ile ayırılabilir.

• Diyelim ki sol alt köşere bir positive example eklendi.

→ Bu anekdot overfitting?  
• Eğer  $C$  çok büyükse ise SVM siyah boundarylerden manevra eder. Ancak bu doğru mu yansıtır. Çünkü eğer  $C$  çok büyükse değilse decision boundary yine siyahken benzer olurdu çünkü term A 0'dan biraz büyük olabilir.



## SVM with Kernels

Hypothesis: Given  $x$ , compute features  $f \in \mathbb{R}^{m+1}$  ( $\theta \in \mathbb{R}^{m+1}$  now!)  
 Predict " $y=1$ " if  $\theta^T \cdot f \geq 0$  where  $\theta^T \cdot f = \theta_0 \cdot f_0 + \theta_1 \cdot f_1 + \dots + \theta_m \cdot f_m$

### Training the Hypothesis:

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

→ name gibb's

• Yukarıdaki ifade genelli inflection: veri ama pratiklik regularization term ekli olur:

•  $\sum_{j=1}^m \theta_j^2 = \theta^T \cdot \theta$  (ignoring  $\theta_0$ )  $\Rightarrow \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$

•  $\sum_{j=1}^m \theta_j^2$  yerine  $\theta^T \cdot M \cdot \theta$  kullanılır yukarıdaki yakını  
 with this modification, it allows to scale to much bigger training sets.  
 Bu sadece matematiksel bir delay.

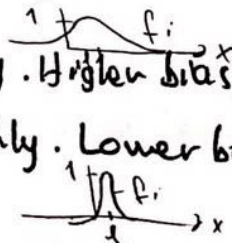
⊗ Can we apply Kernels idea to other algorithms? Like log. reg.

- Uygunabilir. But computational tricks that apply for SVMs don't generalize well to other algorithms like logistic regression. So using kernels with logistic regression will be very slow.
- SVM and Kernels tend to go particularly well together.

⊗ Yukarıdaki cost function'ın minimize etme için off the shelf algorithms var. Bunları kullanmak daha mantıklıdır

### SVM Parameters: Bias vs Variance

- $C$  played a role similar to  $\lambda/x$  where  $\lambda$  is reg. par. for log. reg.  
→ more prone to overfitting
- Large  $C$ : Tend to have a hypothesis with lower bias and higher variance.
- Small  $C$ : → "more" prone to underfitting " higher bias and lower variance
- Large  $\sigma^2$ : Features  $f_i$  vary more smoothly. Higher bias, lower variance.
- Small  $\sigma^2$ : Features  $f_i$  vary less smoothly. Lower bias, higher variance.





## -W7 - Using a SVM -

- What we actually need to do in order to run an SVM.
- It is not recommended to write a software to solve for the parameters  $\Theta$ . Yani  $J(\Theta)$ 'yi minimize edecek algoritmayı bizim yazmamız şart değil. Instead use SVM software package (e.g. liblinear, libsvm, ... ) to solve for parameters  $\Theta$ .

- We still need to specify:
  - Choice of parameter  $C$
  - Choice of kernel (similarity function):

- E.g. No kernel ("linear kernel"): Predict " $y=1$ " if  $\Theta^T \cdot x \geq 0$

- when  $n$  is large,  $m$  is small: Maybe it's better to fit a linear decision boundary and not try to fit a very complicated non-linear function because you might not have enough data, there might be overfitting.

- Another choice is Gaussian Kernel:  $f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$ , where  $l^{(i)} = x^{(i)}$ .

- Need to choose  $\sigma^2$ .

- when  $n$  is small,  $m$  is large: We have enough training set to train a complex nonlinear hypothesis and avoid overfitting. Mesela 2 dimensional training set  $(x_1, x_2)$  ve binlerce training data ran complex bir hipotezi overfit etmeden öğrenilebilir.

- If we decide to use Gaussian kernel here is what we should do:  
function  $f = \text{kernel}(x_1, x_2)$   
 $f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$   
return

Do perform feature scaling before using the Gaussian kernel. Aslında feature'ları normalize etmeliyiz!   
size 1000's → bedrooms → 1-5

- Feature scaling lazım çünkü  $\|x - l\|^2 = (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2$    
Size her scale bayağı olacaktı ki yolda doğru domain eder.

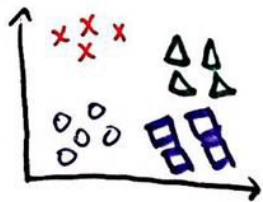
Other choices of kernel: Not all similarity  $(x, l)$  functions make valid kernels. (Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM packages' optimisations run correctly, and do not diverge).

- many off-the-shelf kernels available:

- Polynomial kernel:  $k(x, l) = (x^T \cdot l + \text{constant})^{\text{degree}}$
- More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...



## Multi-class Classification



$y \in \{1, 2, 3, \dots, k\}$

- Many SVM packages already have built-in multi-class classification functionality.
- Otherwise, use one vs all method. (Train  $k$  SVMs, each distinguish different class)
- (Exactly same as Logistic Regression case)

## Logistic Regression vs. SVMs

- $n$  = # of features ( $x \in \mathbb{R}^{n+1}$ ) •  $m$  = # of training examples.

- ⊗ If  $n$  is large (relative to  $m$ ): (E.g.  $n=10000$ ,  $m=10-1000$ )
  - Use logistic regression, or SVM without a kernel ("linear kernel")
  - Because if we have so many features with smaller training sets a linear function will probably do fine. We don't have enough data to fit a complex non-linear function
- ⊗ If  $n$  is small,  $m$  is intermediate ( $n=1-1000$ ,  $m=10-10000$ )
  - Often an SVM with Gaussian kernel will work fine.
- ⊗ If  $n$  is small,  $m$  is large ( $n=1-1000$ ,  $m=50000+$ )
  - In this case SVM with Gaussian kernel would be somehow slow to run. You can go by using SVM with Gaussian kernel for linear.
  - Manually create/add more features, then use logistic regression or SVM without kernel

- ⊗ SVM without kernel is logistic regression got better. Same as SVM with kernel is efficient capabilities

- ⊗ For all of these regimes a Neural Network is likely to work well. But maybe slower to train.

- ⊗ Optimization problem that the SVM has is a convex optimization problem. So a good algorithm will always find the global minimum. In practice local optimas aren't a huge problem for neural networks neither.