

Aprendizado Profundo 1

Overfitting, Underfitting e Regularização

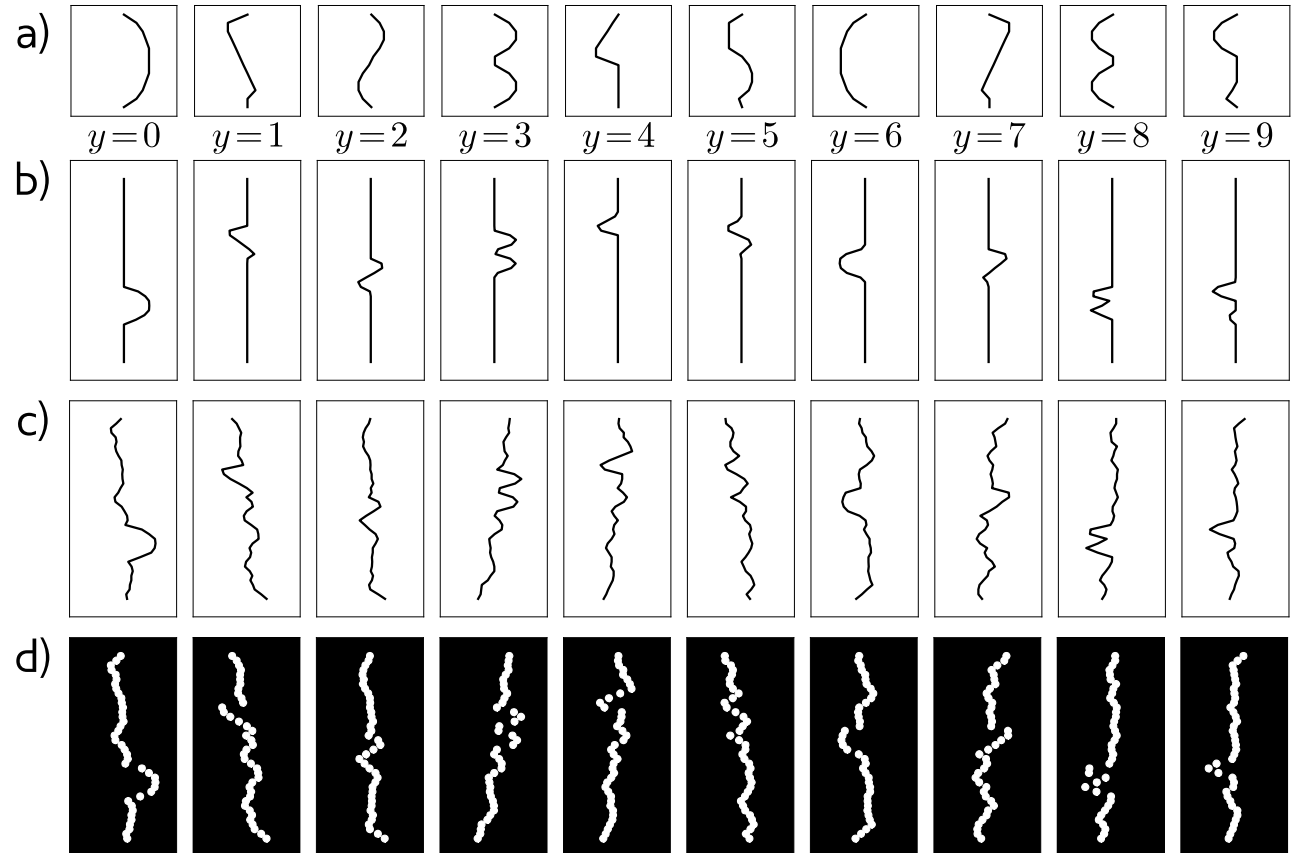
Professor: Lucas Silveira Kupssinskü

Agenda

- Fontes de Erro
 - *Noise, Bias e Variance*
- *Overfitting vs Underfitting*
 - *Bias Variance Tradeoff*
 - *Double Descent*
- Regularização
 - Métodos Explícitos
 - L1, L2
 - Métodos Empíricos
 - *Early Stopping*
 - *Dropout*
 - *Data Augmentation*
 - Métodos Implícitos

Conjunto de dados

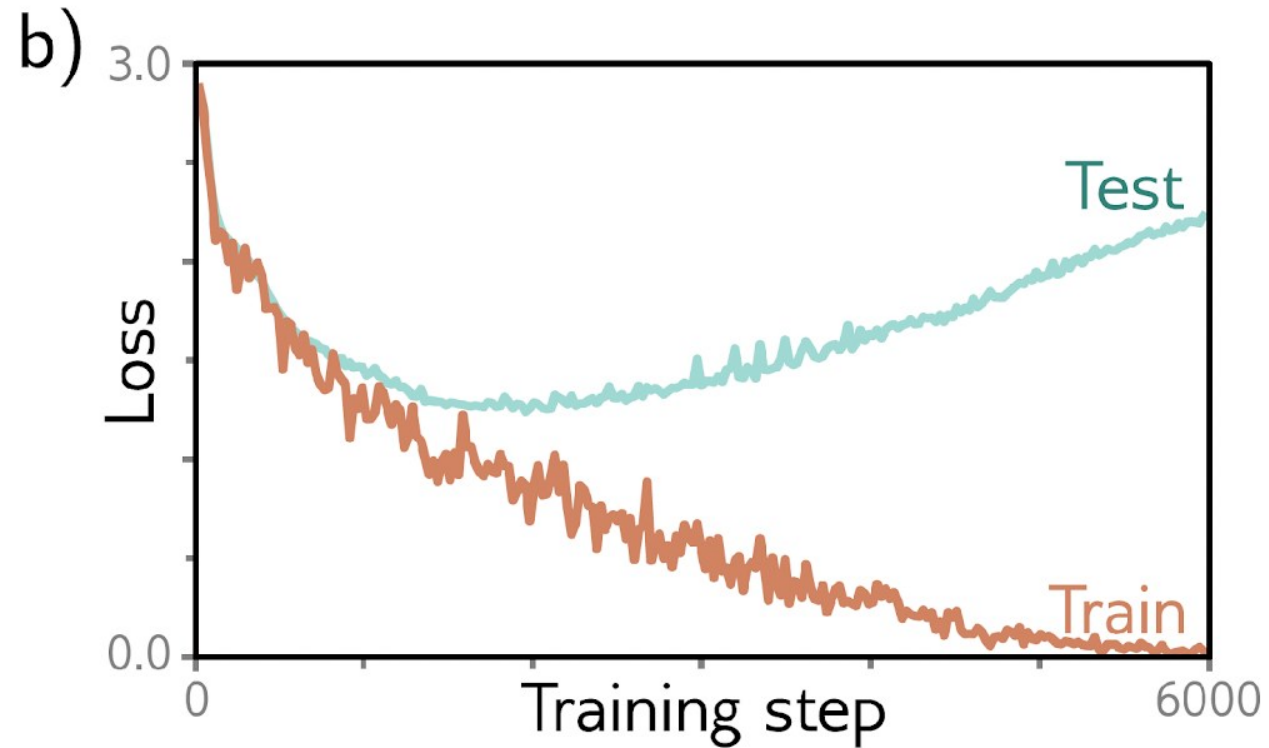
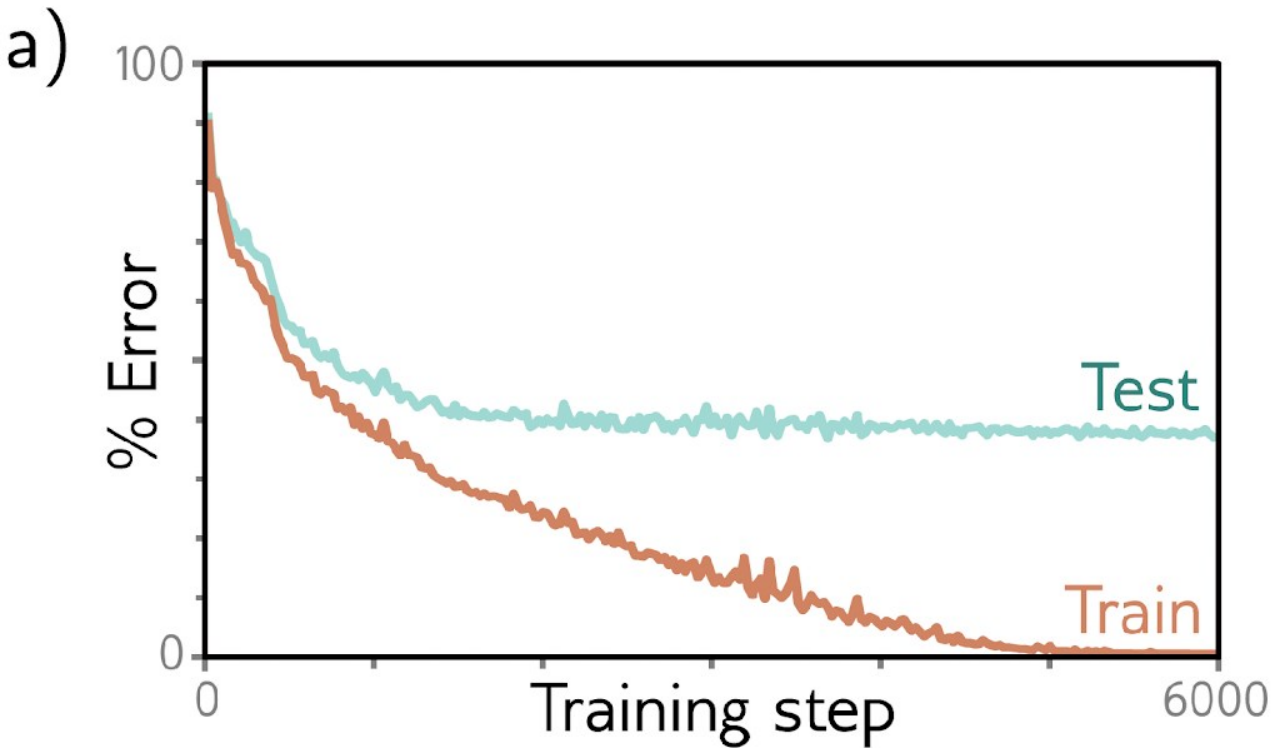
- MNIST 1D
- a) Templates para as 10 classes
- b) Exemplos de treino
- c) Adição de ruído
- d) Amostragem em 40 pontos



MLP

- 40 Entradas
- 10 Saídas (*Softmax*)
- 4000 instâncias de treino
- Duas camadas ocultas
 - 100 unidades em cada
- SGD com *batch size* 100, *learning rate* 0.1
- 6000 iterações
 - Quantas épocas?

Resultados do Treinamento

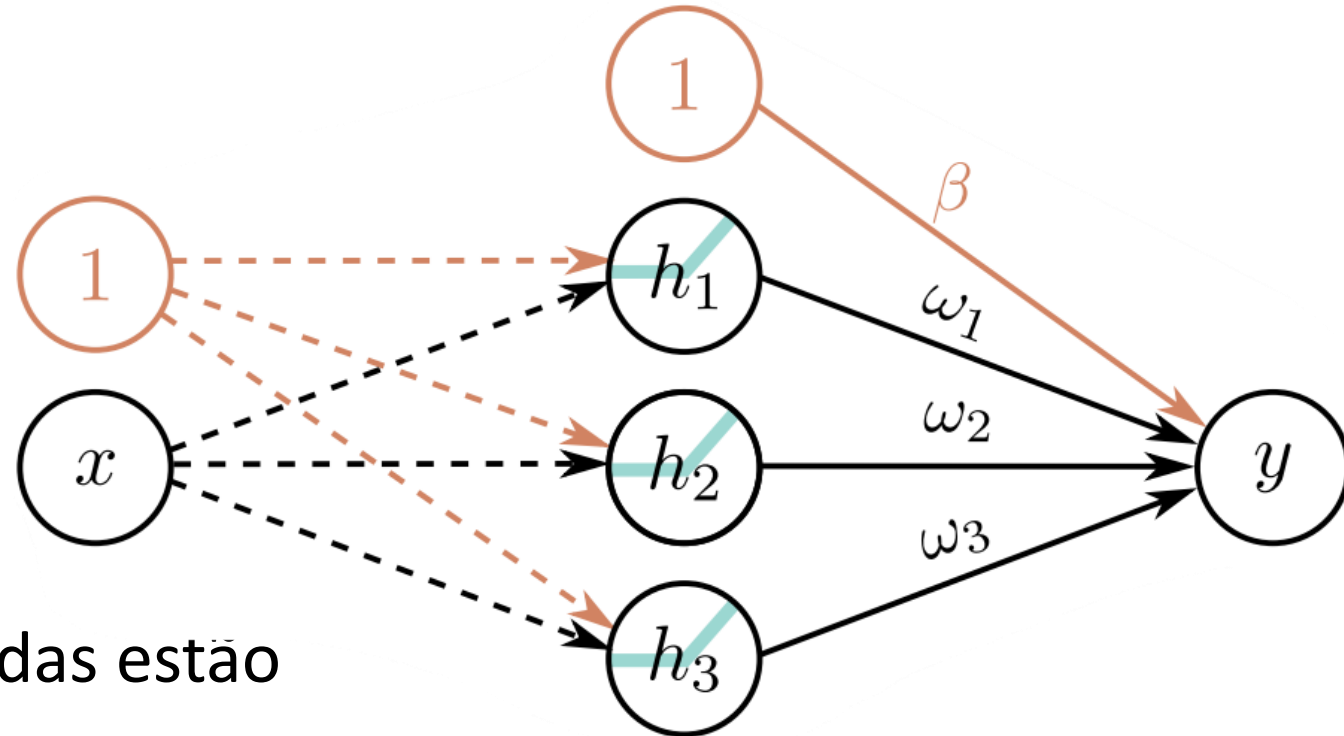


a) A taxa de erro em treinamento diminui até zero

b) A função de custo em teste passa a aumentar enquanto em treinamento ainda há diminuição

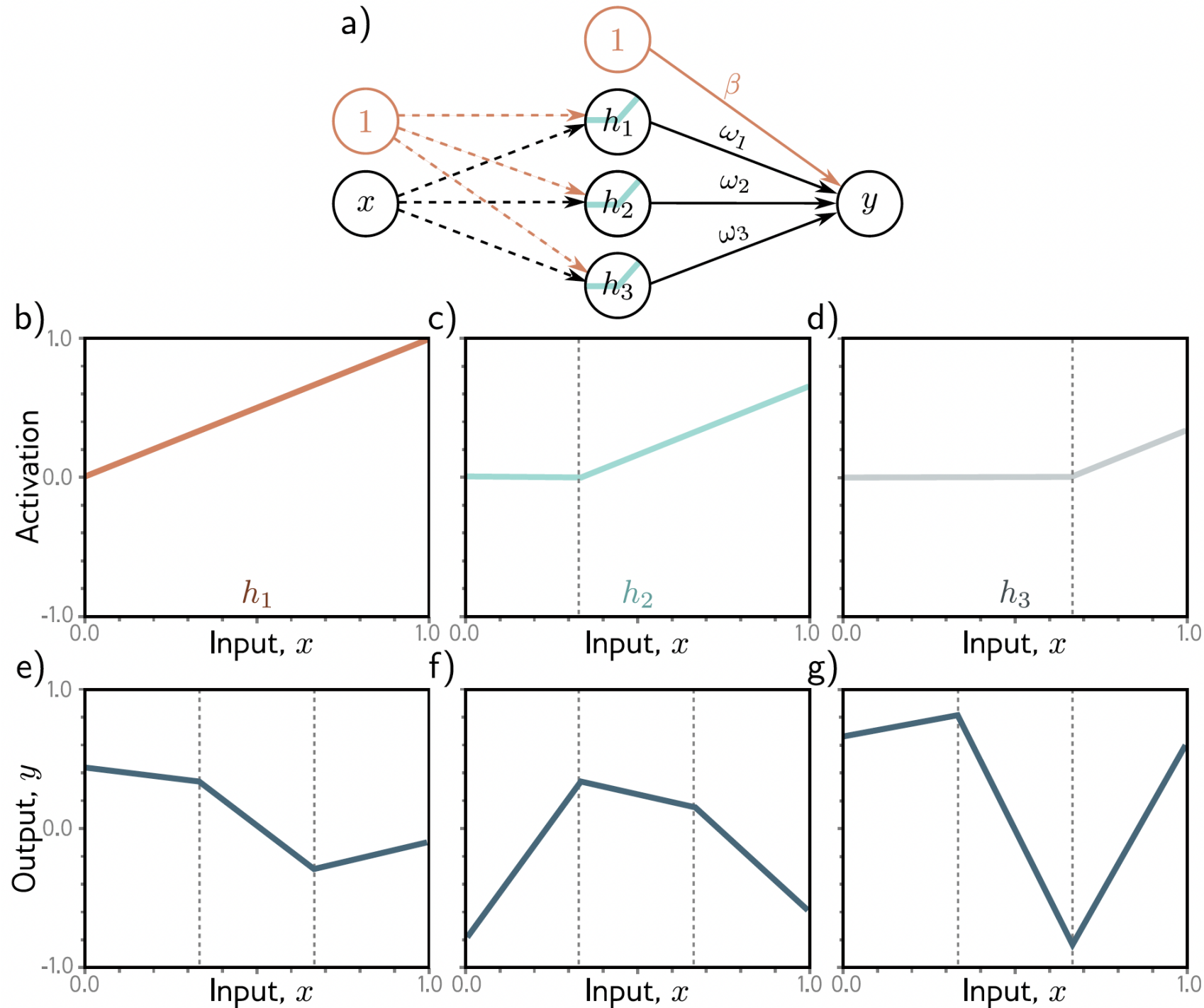
MLP

- 1 Entrada
- 1 Saída (*Linear*)
- Uma camada ocultas
 - 3 unidades em cada
 - ReLU
- Parâmetros nas linhas pontilhadas estão congelados



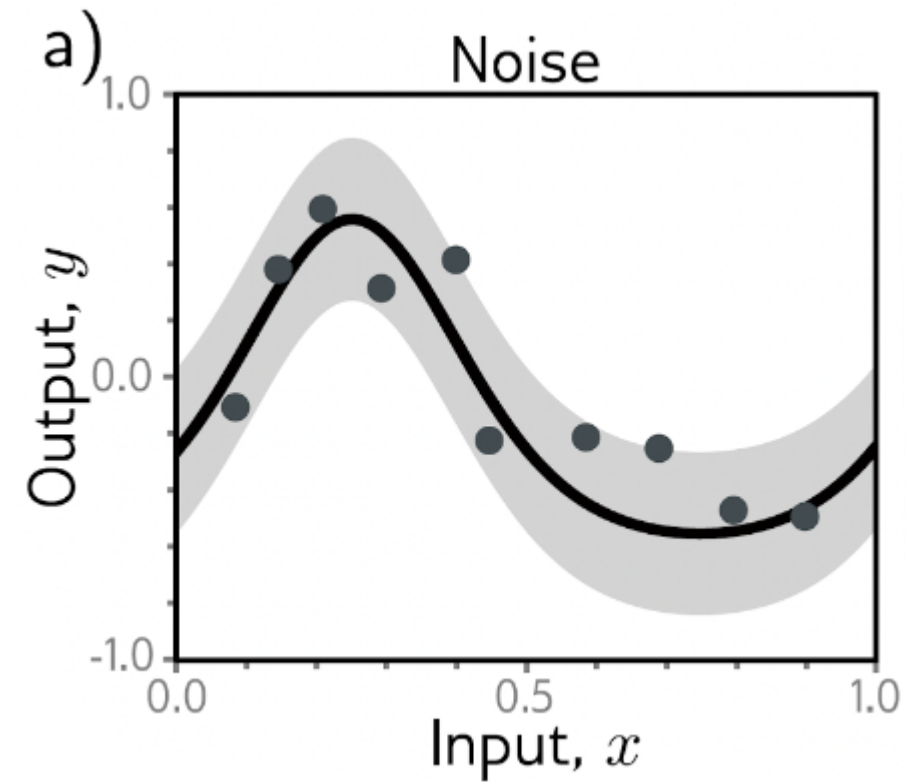
MLP

- Repare em e), f) e g) diferentes parametrizações geram funções lineares diferentes



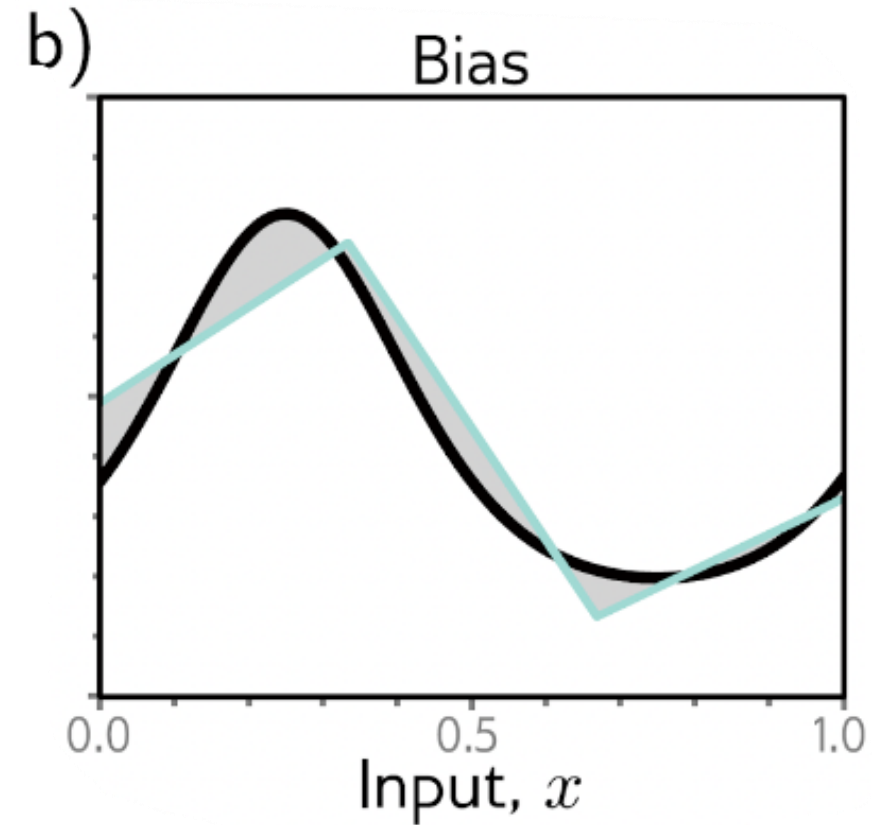
Fontes de Erro

- Ruído
 - Por mais que consigamos modelar a função geradora dos dados (sinusoide em preto), ainda teremos um erro diferente de 0
 - Fatores estocásticos ou ausência de uma variável



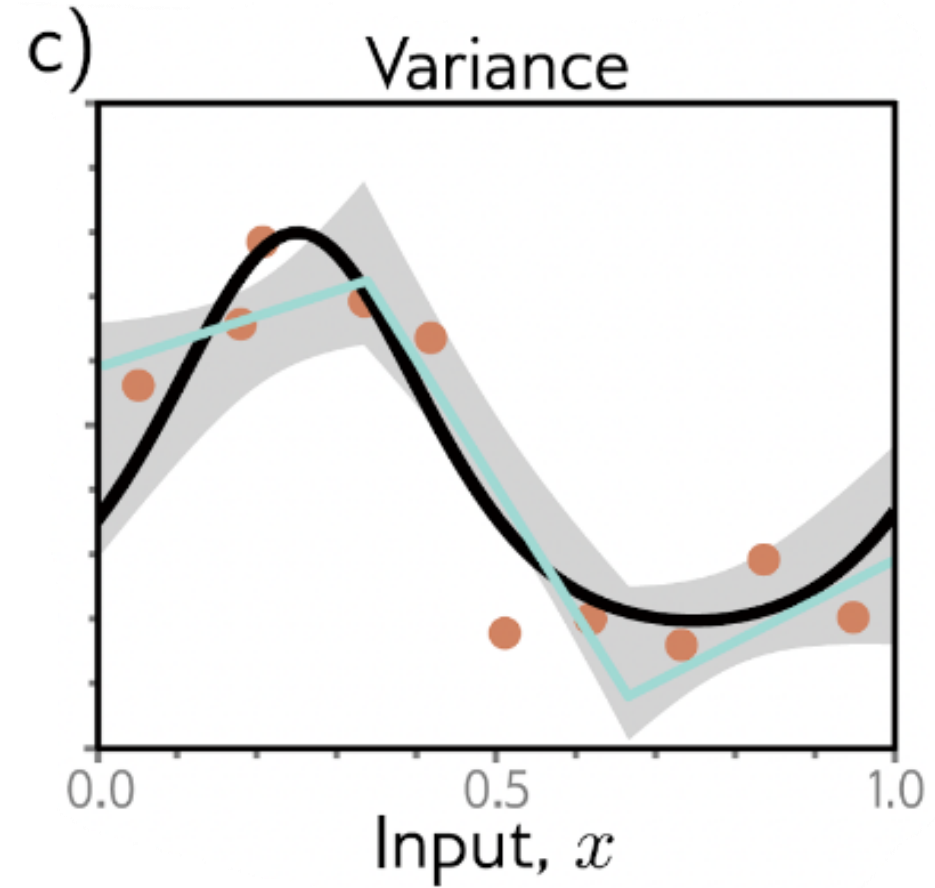
Fontes de Erro

- *Bias*
 - Mesmo com os melhores parâmetros nosso MLP (cyan) não consegue modelar perfeitamente a função em preto
 - Todo modelo precisa de viés indutivo para funcionar

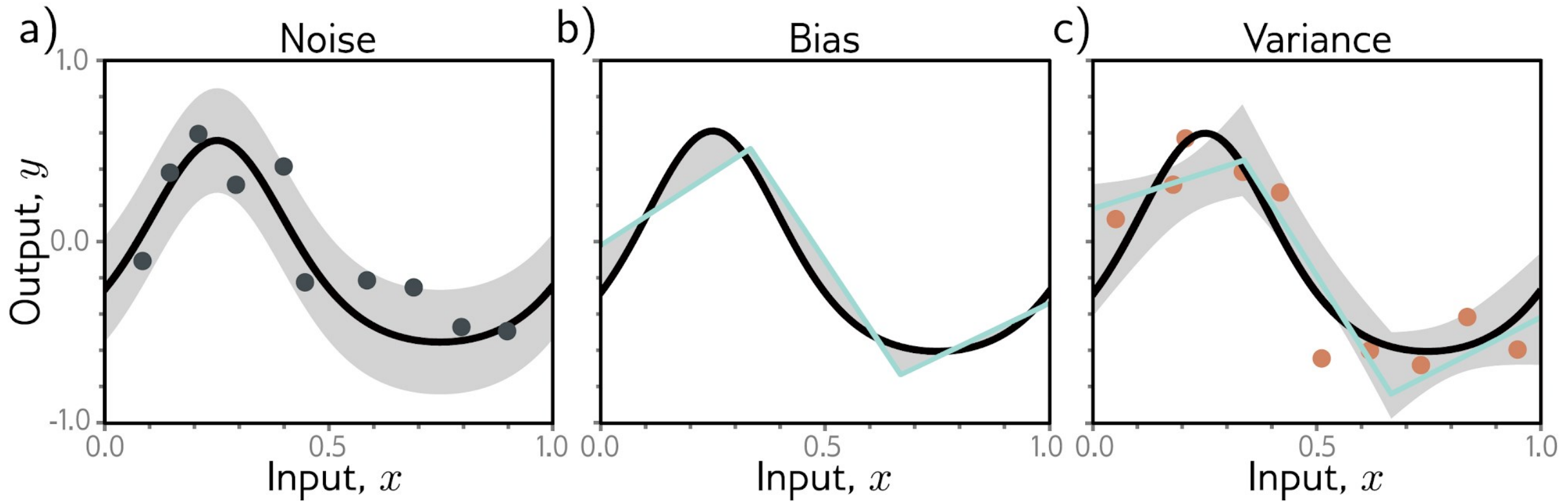


Fontes de Erro

- *Variance*
 - Quando treinamos o modelo em um número limitado de dados, não recuperamos necessariamente o melhor ajuste a função geradora



Fontes de Erro



Fontes de Erro

- Vamos tornar precisa a noção da fonte de erro
- Considere um problema de regressão de uma variável
- Suponha que:
 - Para valores distintos de x temos valores esperados diferentes, ou seja $\mathbb{E}[\text{Pr}(y|x)] = \mu[x]$
 - Existe um ruído fixo $\sigma^2 = \mathbb{E}[(\mu[x] - y[x])^2]$
 - A predição do modelo é denotada por $f[x, \theta]$
 - Vamos considerar o erro quadrático como função de custo $L[x] = (f[x, \theta] - y[x])^2$

Fontes de Erro

- Suponha que:
 - Para valores distintos de x temos valores esperados diferentes, ou seja $\mathbb{E}[\text{Pr}(y|x)] = \mu[x]$
 - Existe um ruído fixo $\sigma^2 = \mathbb{E}[(\mu[x] - y[x])^2]$
 - A predição do modelo é denotada por $f[x, \theta]$
 - Vamos considerar o erro quadrático como função de custo $L[x] = (f[x, \theta] - y[x])^2$

$$L[x] = (f[x, \theta] - y[x])^2$$

Fontes de Erro

- Suponha que:
 - Para valores distintos de x temos valores esperados diferentes, ou seja $\mathbb{E}[\Pr(y|x)] = \mu[x]$
 - Existe um ruído fixo $\sigma^2 = \mathbb{E}[(\mu[x] - y[x])^2]$
 - A predição do modelo é denotada por $f[x, \theta]$
 - Vamos considerar o erro quadrático como função de custo $L[x] = (f[x, \theta] - y[x])^2$

$$\begin{aligned} L[x] &= (f[x, \theta] - y[x])^2 \\ &= ((f[x, \theta] - \mu[x]) + (\mu[x] - y[x]))^2 \end{aligned}$$

Fontes de Erro

- Suponha que:
 - Para valores distintos de x temos valores esperados diferentes, ou seja $\mathbb{E}[\text{Pr}(y|x)] = \mu[x]$
 - Existe um ruído fixo $\sigma^2 = \mathbb{E}[(\mu[x] - y[x])^2]$
 - A predição do modelo é denotada por $f[x, \theta]$
 - Vamos considerar o erro quadrático como função de custo $L[x] = (f[x, \theta] - y[x])^2$

$$\begin{aligned} L[x] &= (f[x, \theta] - y[x])^2 \\ &= ((f[x, \theta] - \mu[x]) + (\mu[x] - y[x]))^2 \\ &= (f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2 \end{aligned}$$

Fontes de Erro

- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y

$$\begin{aligned} L[x] &= (f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2 \\ \mathbb{E}_y[L[x]] &= \mathbb{E}_y[(f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= \mathbb{E}_y[(f[x, \theta] - \mu[x])^2] + \mathbb{E}_y[2(f[x, \theta] - \mu[x])(\mu[x] - y[x])] + \mathbb{E}_y[(\mu[x] - y[x])^2] \end{aligned}$$

Não dependem de y

Fontes de Erro

- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y

$$\begin{aligned} L[x] &= (f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2 \\ \mathbb{E}_y[L[x]] &= \mathbb{E}_y[(f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= \mathbb{E}_y[(f[x, \theta] - \mu[x])^2] + \mathbb{E}_y[2(f[x, \theta] - \mu[x])(\mu[x] - y[x])] + \mathbb{E}_y[(\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= (f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - \cancel{\mathbb{E}_y[y[x]]}) + \mathbb{E}_y[(\mu[x] - y[x])^2] \\ &\quad \mu[x] \end{aligned}$$

Fontes de Erro

- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y

$$\begin{aligned} L[x] &= (f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2 \\ \mathbb{E}_y[L[x]] &= \mathbb{E}_y[(f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= \mathbb{E}_y[(f[x, \theta] - \mu[x])^2] + \mathbb{E}_y[2(f[x, \theta] - \mu[x])(\mu[x] - y[x])] + \mathbb{E}_y[(\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= (f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - \mathbb{E}_y[y[x]]) + \mathbb{E}_y[(\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= (f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])\mathbf{0} + \cancel{\mathbb{E}_y[(\mu[x] - y[x])^2]} \sigma^2 \end{aligned}$$

Fontes de Erro

- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y

$$\begin{aligned} L[x] &= ((f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2) \\ \mathbb{E}_y[L[x]] &= \mathbb{E}_y[(f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= \mathbb{E}_y[(f[x, \theta] - \mu[x])^2] + \mathbb{E}_y[2(f[x, \theta] - \mu[x])(\mu[x] - y[x])] + \mathbb{E}_y[(\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= (f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])(\mu[x] - \mathbb{E}_y[y[x]]) + \mathbb{E}_y[(\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= (f[x, \theta] - \mu[x])^2 + 2(f[x, \theta] - \mu[x])0 + \mathbb{E}_y[(\mu[x] - y[x])^2] \\ \mathbb{E}_y[L[x]] &= (f[x, \theta] - \mu[x])^2 + \sigma^2 \end{aligned}$$

Este termo se refere a nossas previsões.
Aqui podemos trabalhar mais

Este termo se refere ao ruído. Não
temos como diminuir essa componente

Fontes de Erro

- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y
 - Vamos definir um modelo médio em relação a todos os possíveis *datasets* $f_\mu[x] = \mathbb{E}_D [f[x, \theta[D]]]$

$$\begin{aligned}\mathbb{E}_y[L[x]] &= (f[x, \theta] - \mu[x])^2 + \sigma^2 \\ (f[x, \theta] - \mu[x])^2 &= \left((f[x, \theta] - f_\mu[x]) + (f_\mu[x] - \mu[x]) \right)^2 \\ &= (f[x, \theta] - f_\mu[x])^2 + 2(f[x, \theta] - f_\mu[x])(f_\mu[x] - \mu[x]) + (f_\mu[x] - \mu[x])^2\end{aligned}$$

Fontes de Erro

- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y
 - Vamos definir um modelo médio em relação a todos os possíveis *datasets* $f_\mu[x] = \mathbb{E}_D [f[x, \theta[D]]]$

$$\begin{aligned}\mathbb{E}_y[L[x]] &= (f[x, \theta] - \mu[x])^2 + \sigma^2 \\ (f[x, \theta] - \mu[x])^2 &= \left((f[x, \theta] - f_\mu[x]) + (f_\mu[x] - \mu[x]) \right)^2 \\ &= (f[x, \theta] - f_\mu[x])^2 + 2(f[x, \theta] - f_\mu[x])(f_\mu[x] - \mu[x]) + (f_\mu[x] - \mu[x])^2 \\ \mathbb{E}_D[(f[x, \theta] - \mu[x])^2] &= \mathbb{E}_D \left[(f[x, \theta] - f_\mu[x])^2 + 2(f[x, \theta] - f_\mu[x])(f_\mu[x] - \mu[x]) + (f_\mu[x] - \mu[x])^2 \right] \\ &= \mathbb{E}_D \left[(f[x, \theta] - f_\mu[x])^2 \right] + \mathbb{E}_D \left[2(f[x, \theta] - f_\mu[x])(f_\mu[x] - \mu[x]) \right] + \mathbb{E}_D \left[(f_\mu[x] - \mu[x])^2 \right] \\ \mathbb{E}_D[(f[x, \theta] - \mu[x])^2] &= \mathbb{E}_D \left[(f[x, \theta] - f_\mu[x])^2 \right] + (f_\mu[x] - \mu[x])^2\end{aligned}$$

Fontes de Erro

- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y
 - Vamos definir um modelo médio em relação a todos os possíveis *datasets* $f_\mu[x] = \mathbb{E}_D [f[x, \theta[D]]]$

$$\mathbb{E}_D [\mathbb{E}_y [L[x]]] = \mathbb{E}_D [(f[x, \theta] - f_\mu[x])^2] + (f_\mu[x] - \mu[x])^2 + \sigma^2$$

Quanto o modelo atual varia
em relação ao modelo médio

Variância

Fontes de Erro

- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y
 - Vamos definir um modelo médio em relação a todos os possíveis *datasets* $f_\mu[x] = \mathbb{E}_D [f[x, \theta[D]]]$

$$\mathbb{E}_D [\mathbb{E}_y [L[x]]] = \mathbb{E}_D [(f[x, \theta] - f_\mu[x])^2] + \boxed{(f_\mu[x] - \mu[x])^2} + \sigma^2$$

Quanto o modelo médio difere da média real da variável alvo

Viés

Fontes de Erro

- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y
 - Vamos definir um modelo médio em relação a todos os possíveis *datasets* $f_\mu[x] = \mathbb{E}_D [f[x, \theta[D]]]$

$$\mathbb{E}_D [\mathbb{E}_y [L[x]]] = \mathbb{E}_D [(f[x, \theta] - f_\mu[x])^2] + (f_\mu[x] - \mu[x])^2 + \sigma^2$$

Quanto o modelo médio difere da média real da variável alvo

ruído

Fontes de Erro

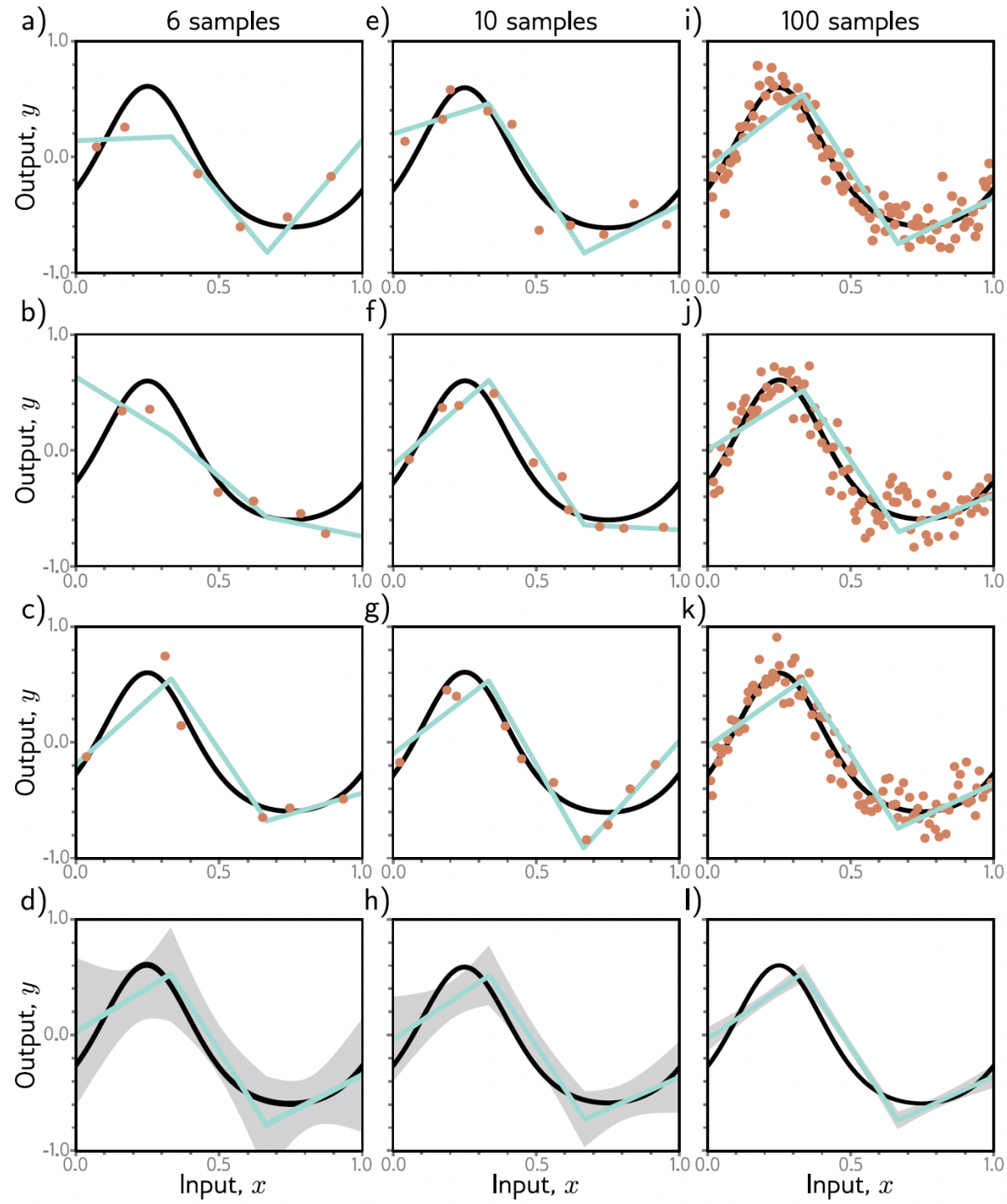
- Agora vamos verificar qual é o valor esperado da função de custo em relação a um valor y

$$\mathbb{E}_D \left[\mathbb{E}_y[L[x]] \right] = \underbrace{\mathbb{E}_D \left[(f[x, \theta] - f_\mu[x])^2 \right]}_{\text{Variância}} + \underbrace{(f_\mu[x] - \mu[x])^2}_{\text{Viés}} + \underbrace{\sigma^2}_{\text{ruído}}$$

Em modelos lineares essas fontes de erro são três componentes aditivos, em modelos não lineares essa interação pode ser mais complexa

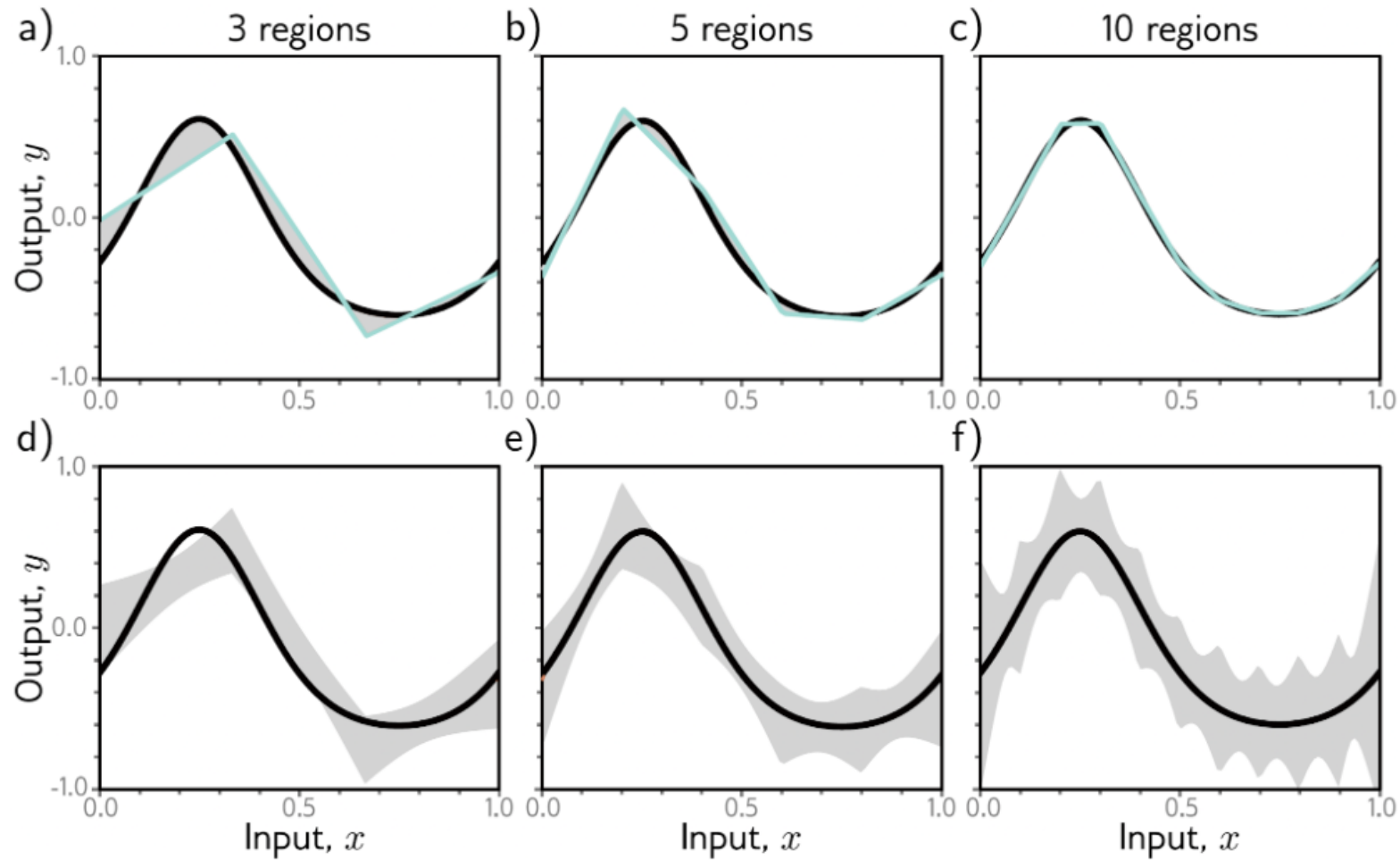
Dataset Size vs Variance

- Lembre-se que estamos falando do modelo que tem capacidade de modelar 3 regiões lineares
- Grande variância para 6 instâncias
- Pequena variância para 100 instâncias



Capacidade do Modelo vs Variância

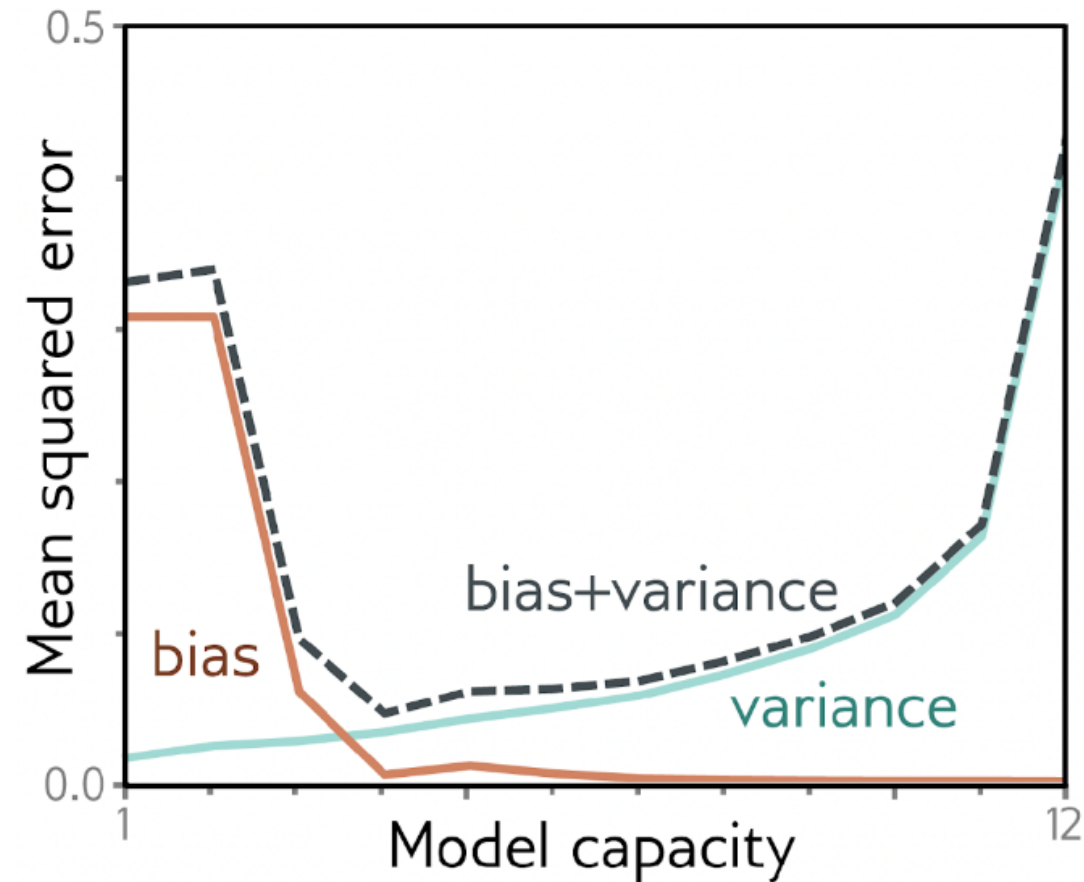
- + capacidade
- + variância



Trade-off Viés-Variância

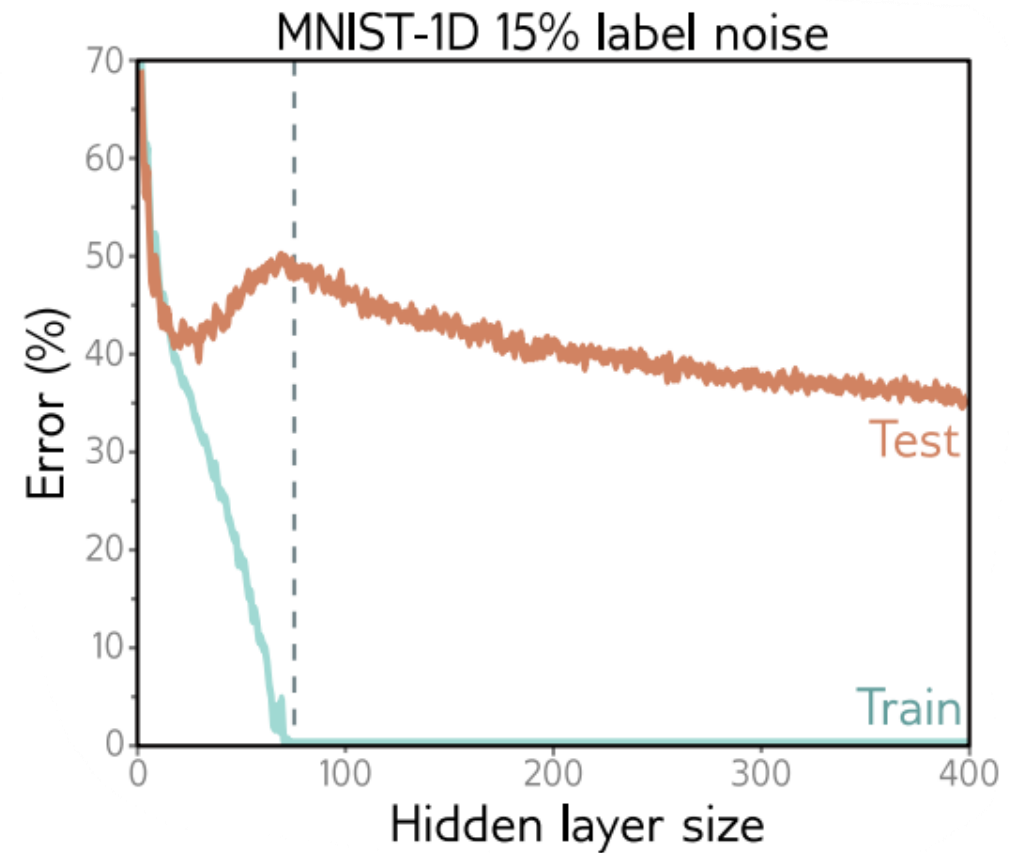
- **Usualmente**

- Quanto maior a capacidade, maior o erro por variância
- Quanto menor a capacidade, maior o erro por viés



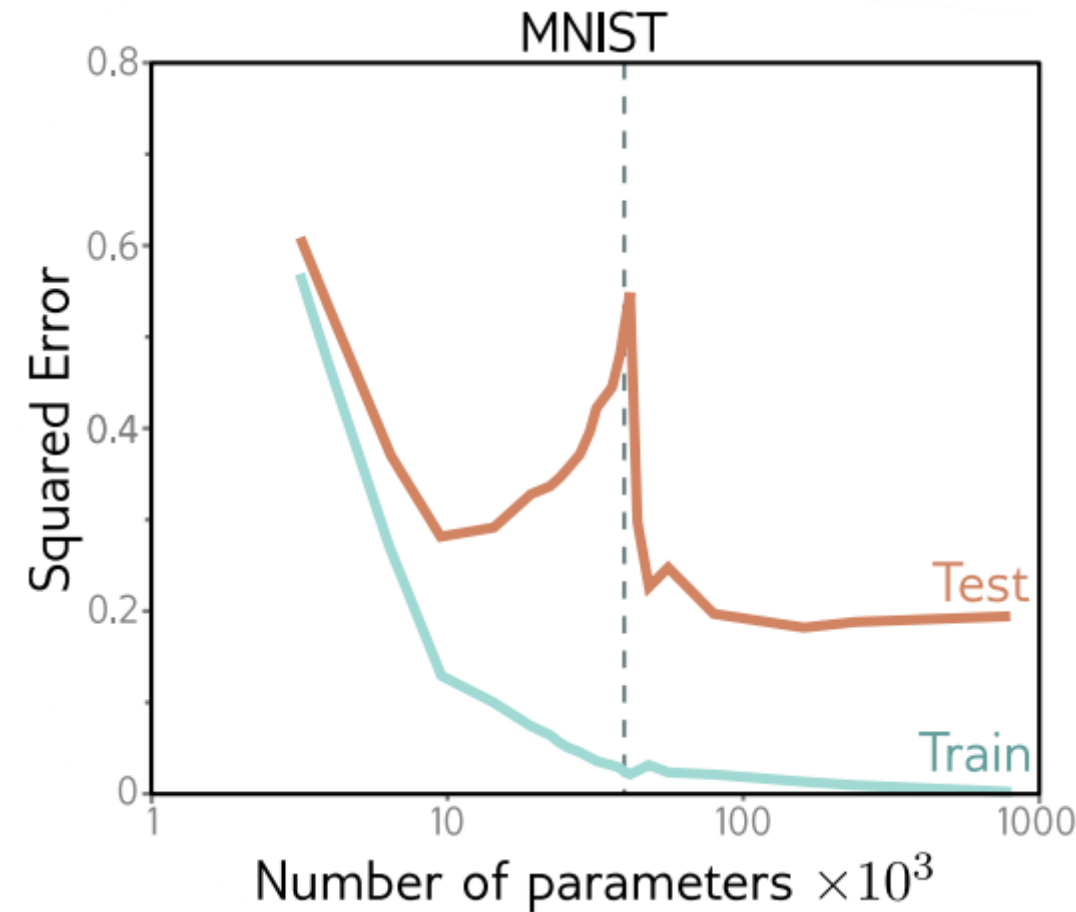
Double Descent

- Será que temos esse trade-off com redes neurais?
 - Observe o comportamento no gráfico
 - MNIST-1D



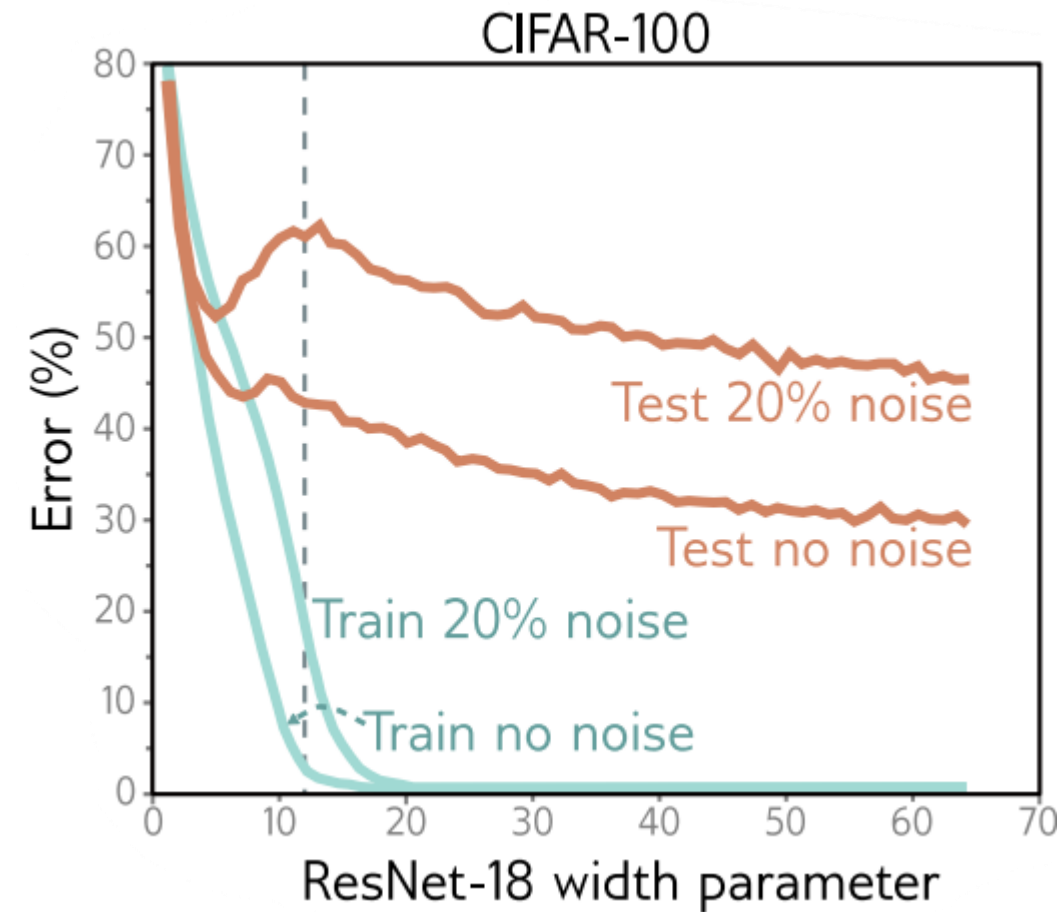
Double Descent

- Será que temos esse trade-off com redes neurais?
 - Observe o comportamento no gráfico
 - MNIST-1D
 - MNIST



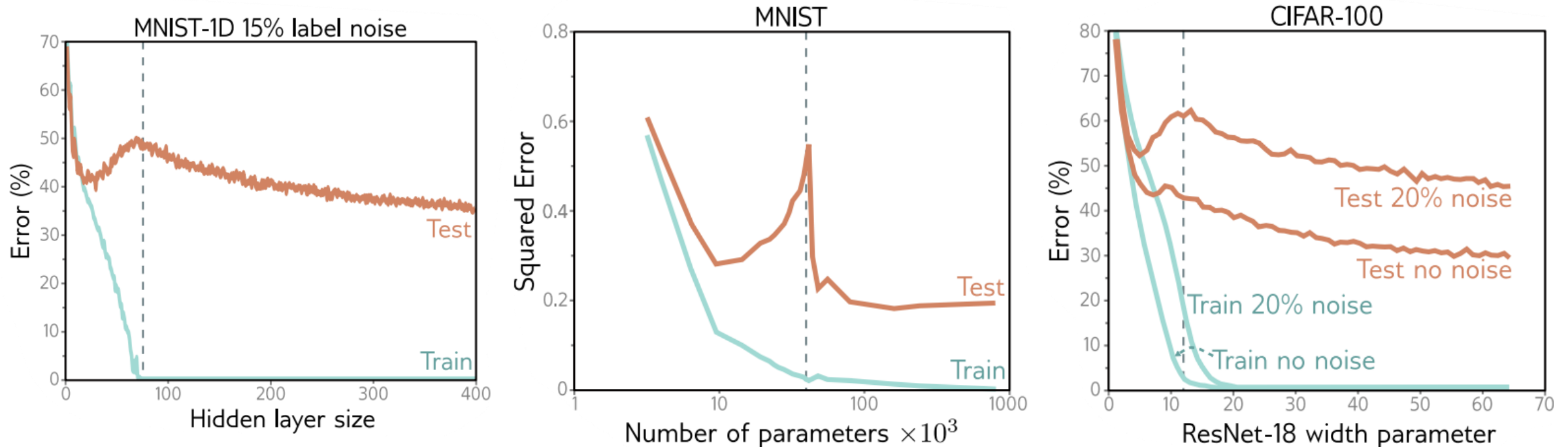
Double Descent

- Será que temos esse trade-off com redes neurais?
 - Observe o comportamento no gráfico
 - MNIST-1D
 - MNIST
 - CIFAR-100



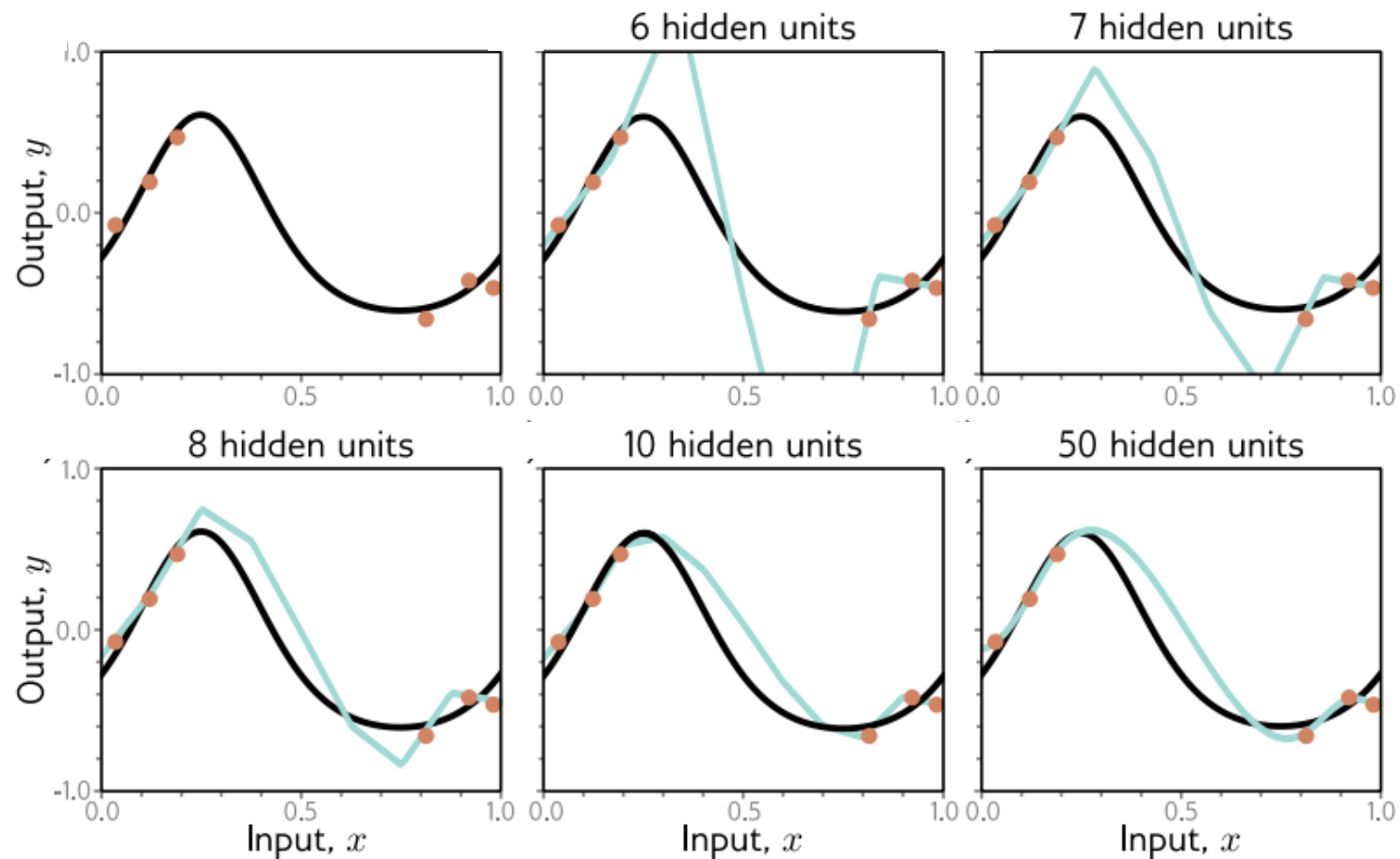
Double Descent

- A linha pontilhada marca o local onde o número de parâmetros da rede é igual ao número de instâncias de treino



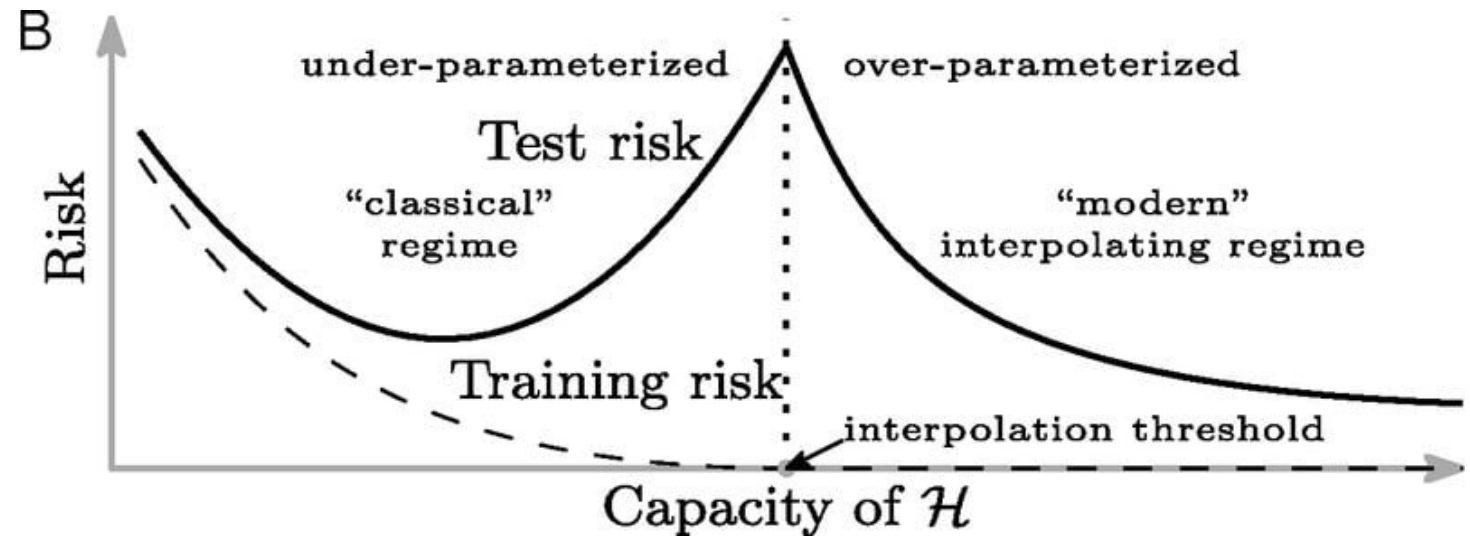
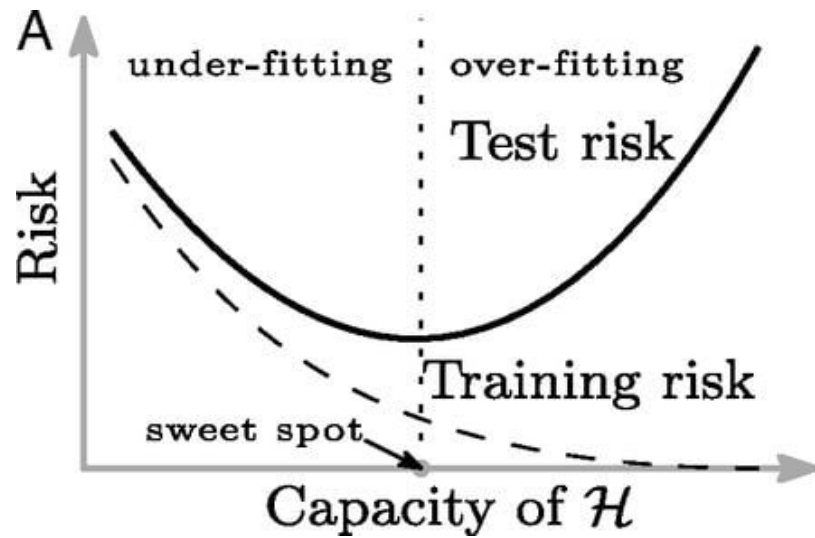
Double Descent

- Um outro experimento para ilustrar



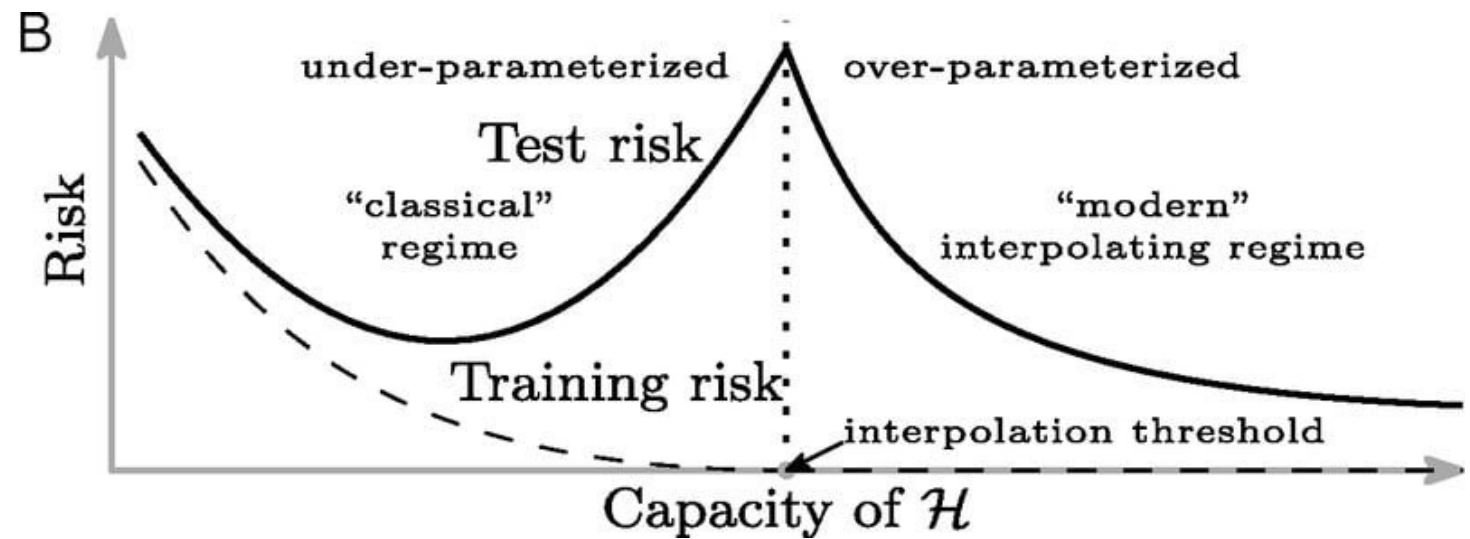
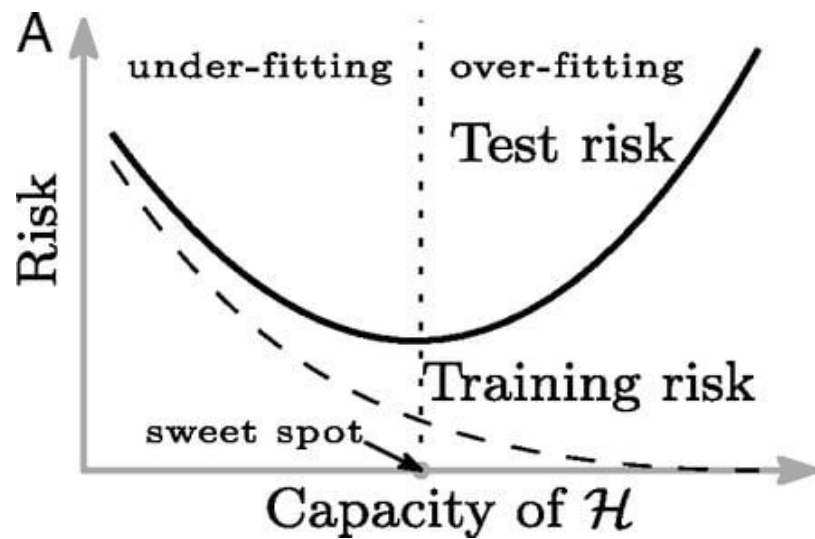
Double Descent

- Por que esse fenômeno ocorre?



Double Descent

- Talvez a inicialização encoraje funções mais suaves
- Talvez o treinamento encoraje funções mais suaves
- Ninguém sabe ao certo, são apenas bons palpites 😊



Regularização

- L1 e L2
 - Enviamos nossa função de custo para dar preferência para parâmetros de baixa magnitude
 - Do ponto de vista probabilístico, estamos adicionando um *prior* no critério de máxima verossimilhança
 - Como redes neurais são atuam em problemas gerais, nosso *prior* tem que ser algo “genérico”

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left[\sum_{i=1}^N l(f(x^{(i)}, \theta), y^{(i)}) + \lambda \|\theta\|^2 \right]$$

$$\theta_{t+1} = \theta - \eta \frac{\partial(\sum_{i=1}^N l(f(x^{(i)}, \theta), y^{(i)}) + \lambda \|\theta\|^2)}{\partial \theta}$$

$$\theta_{t+1} = \theta - \eta \left(\frac{\partial(\sum_{i=1}^N l(f(x^{(i)}, \theta), y^{(i)}))}{\partial \theta} + \frac{\partial(\lambda \|\theta\|^2)}{\partial \theta} \right)$$

$$\theta_{t+1} = \theta - \eta \left(\frac{\partial(\sum_{i=1}^N l(f(x^{(i)}, \theta), y^{(i)}))}{\partial \theta} + 2\lambda\theta \right)$$

$$\theta_{t+1} = \theta - \eta \frac{\partial(\sum_{i=1}^N l(f(x^{(i)}, \theta), y^{(i)}))}{\partial \theta} - \eta 2\lambda\theta$$

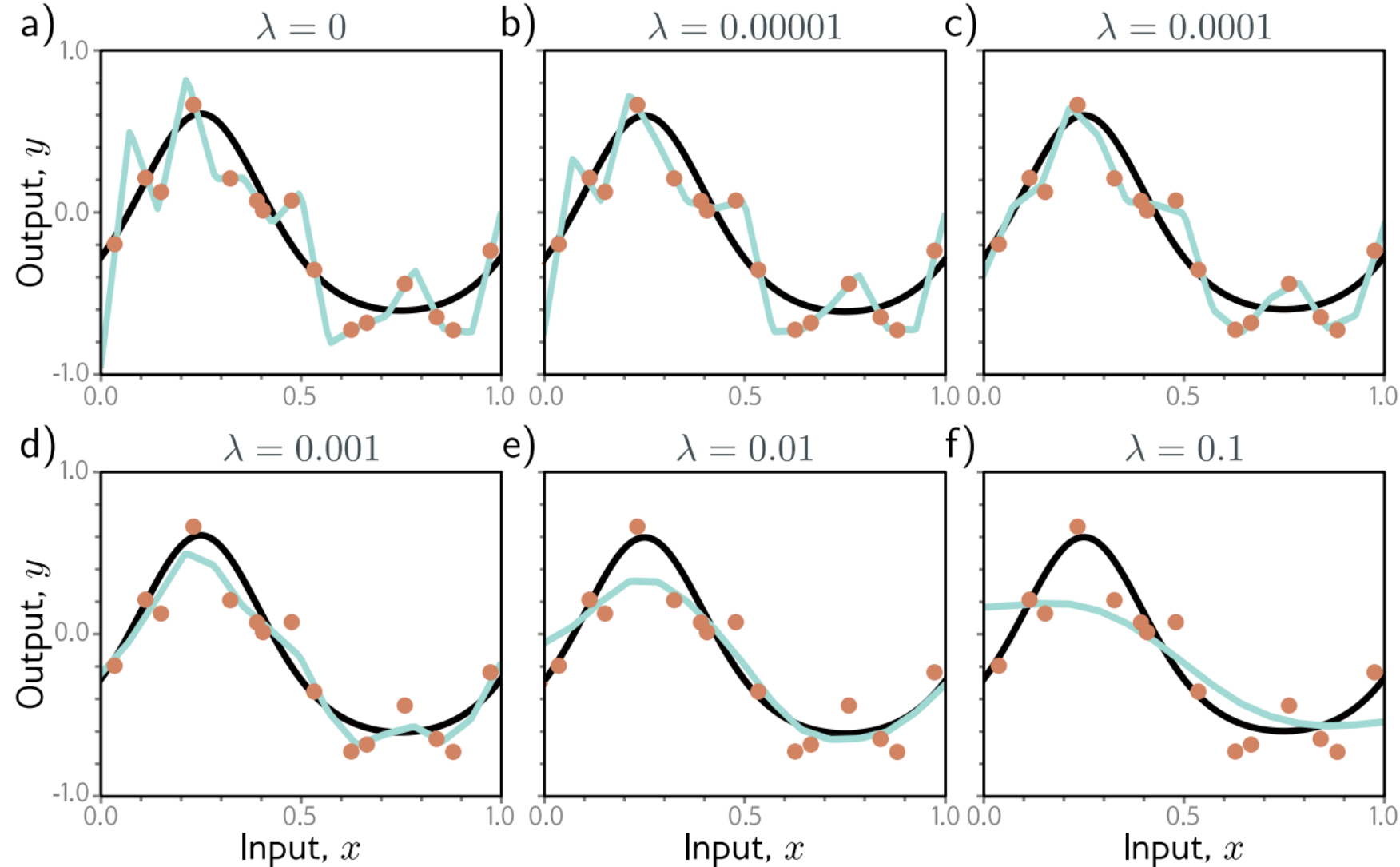
$$\theta_{t+1} = \theta(1 - 2\eta\lambda) - \eta \frac{\partial(\sum_{i=1}^N l(f(x^{(i)}, \theta), y^{(i)}))}{\partial \theta}$$

$$\theta_{t+1} = \theta(1 - 2\eta\lambda) - \eta \nabla J$$

Regularização

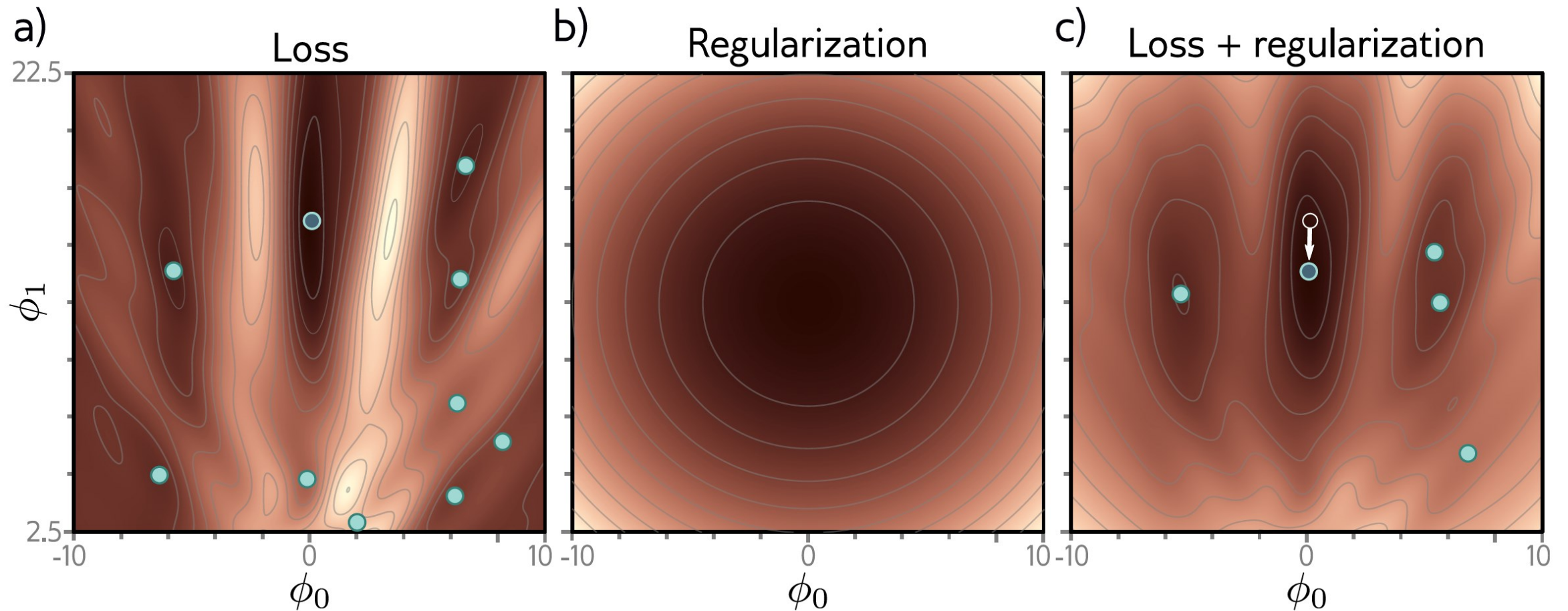
- L1 e L2

- Perceba ao lado o efeito de aumentar o parâmetro de regularização



Regularização

- L1 e L2

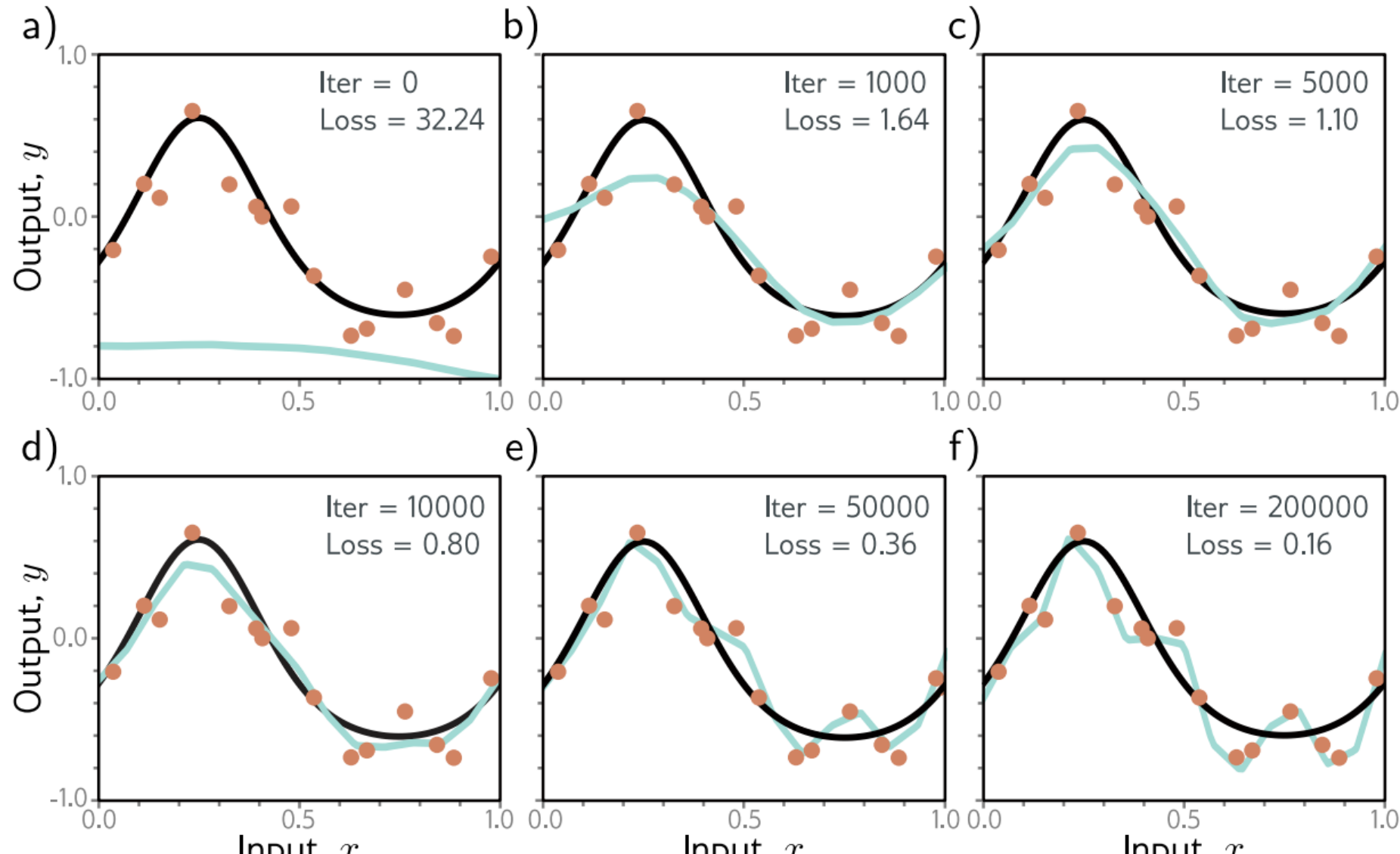


Regularização

- *Early Stopping*
 - Consiste em monitorar o desempenho da rede neural em um conjunto de dados de validação e parar o treinamento antes da loss “convergir”
 - Tem um efeito similar a regularização L2, pois não deixa os parâmetros crescerem indiscriminadamente
 - Usualmente possui apenas um hiperparâmetro associado, que define quantas iterações vamos aguardar sem que haja melhoria na função de custo em validação

Regularização

- *Early Stopping*

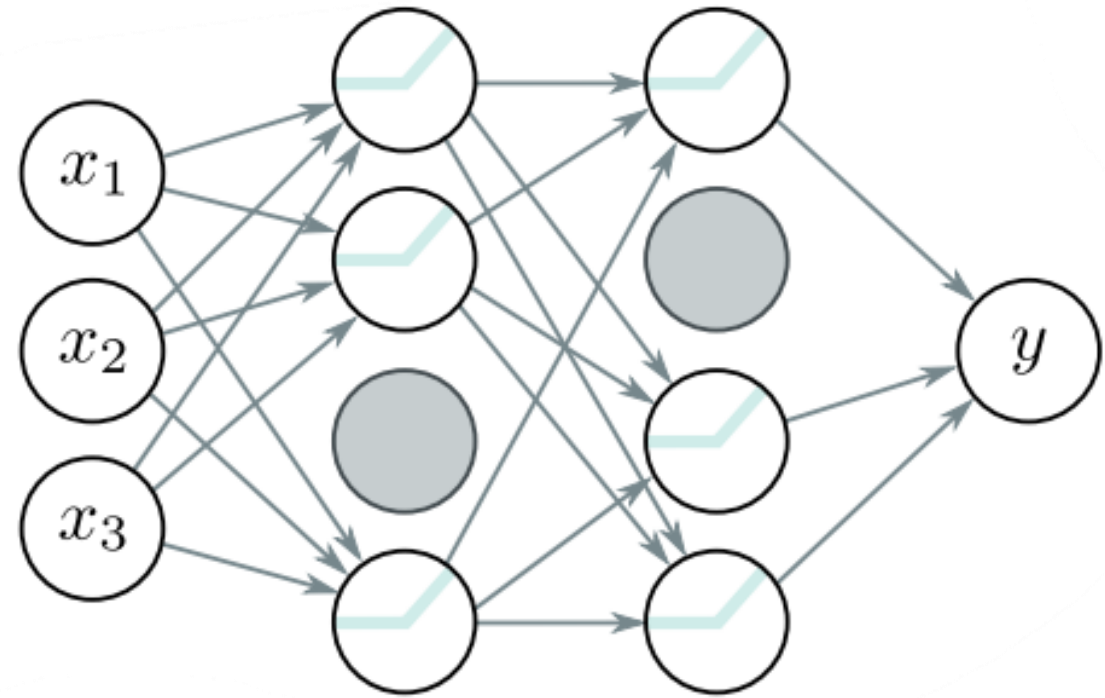
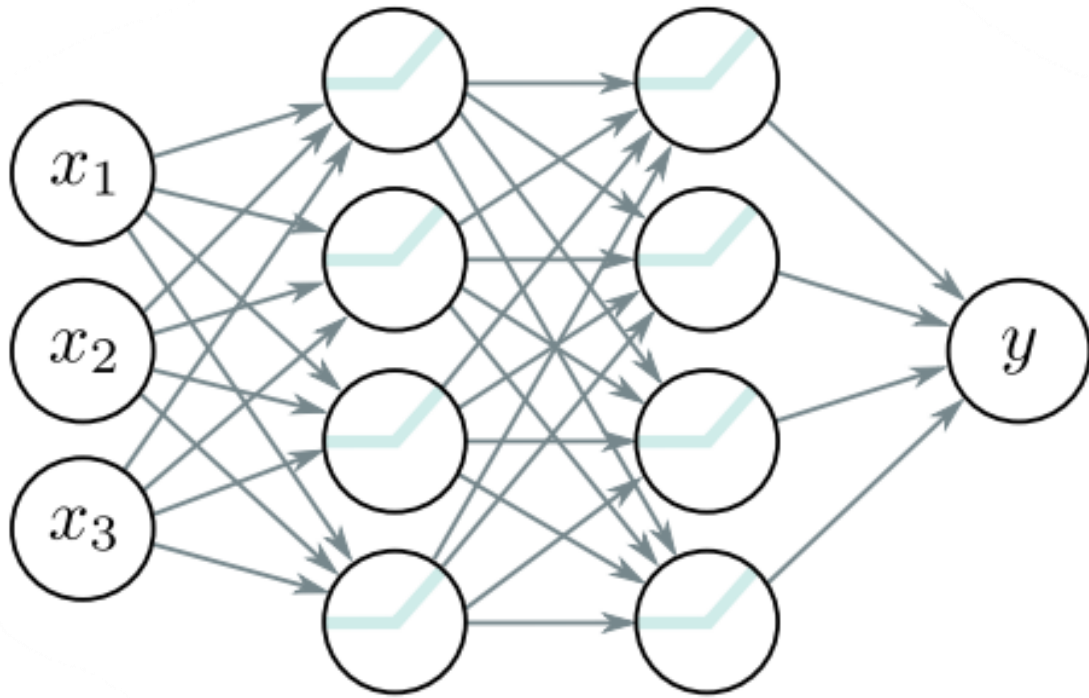


Regularização

- *Drop Out*
 - Durante o treinamento, são desligadas (ou desconectadas) de forma aleatória algumas unidades em uma camada oculta
 - Tem o efeito prático de multiplicarmos as ativações por um ruído que segue uma distribuição de Bernoulli

Regularização

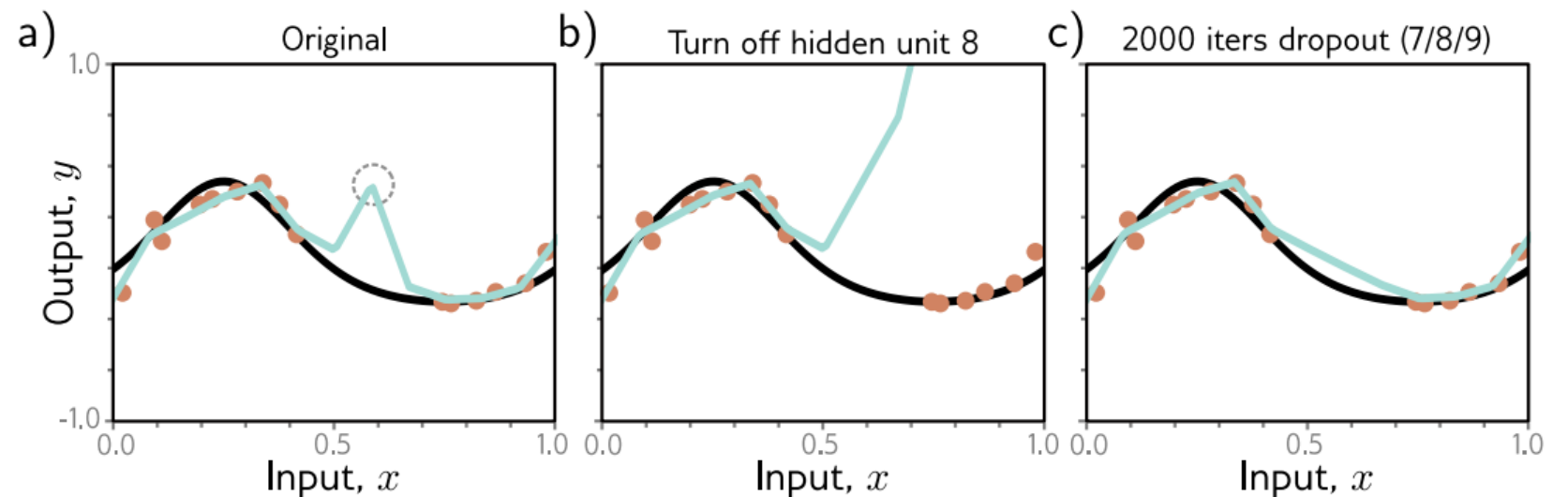
- *Drop Out*



Regularização

- *Drop Out*

- Perceba que na figura a) abaixo, o modelo já se ajusta perfeitamente aos dados de treinamento (continuar treinando não deve ajudar)
- Contudo, se uma das unidades for desligada, o pico no meio da função causará um grande aumento na *loss* (ver b)
- Se continuarmos desligando aleatoriamente as unidades, esse pico no meio da função vai acabar desaparecendo (ver c)



Regularização

- *Drop Out*
 - Na inferência temos duas opções para compensar o aumento da intensidade das ativações:
 - Multiplicamos as ativações por $1 - dropout_rate$
 - Usar *Monte Carlo Drop Out*, rodar o *forward* da rede múltiplas vezes com diferentes unidades desligadas e combinar as diferentes saídas (de forma análoga a um *ensemble*)

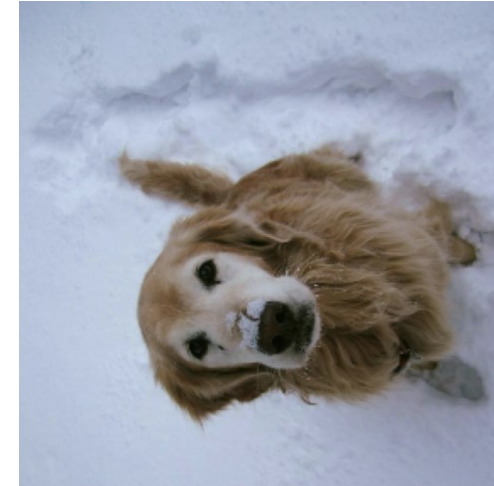
Regularização

- *Data Augmentation*

- Consiste em gerar novas instâncias fazendo manipulações nos dados
 - Rotações, Flip, Mudança de Cor, Equalização, Random Crop, Adição de Ruído,...
 - Adicionamos as invariâncias desejadas via augmentation



Dog



Dog

Regularização

- Implícita
 - É um fato peculiar que a descida de gradiente encontre soluções que generalizam para além do conjunto de treinamento
 - Observou-se que tanto a descida de gradiente em batch (N) e a descida de gradiente estocástica (1) não fazem uma trajetória direta ao mínimo local da função, existe uma preferência para algumas soluções em detrimento de outras
 - Essas preferências recebem o nome de regularização implícita

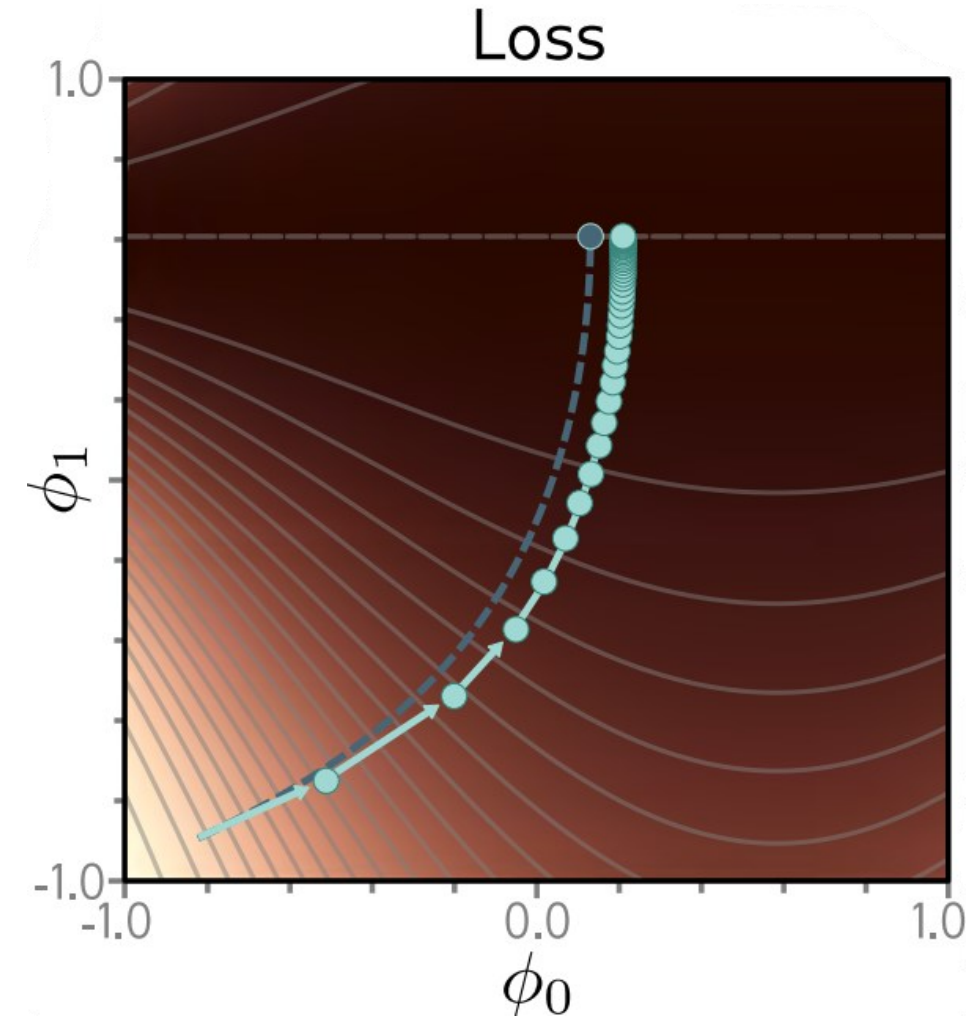
Regularização Implícita

- A descida de gradiente contínua pode ser representada a partir da seguinte equação diferencial

$$\frac{\partial \theta}{\partial t} = - \frac{\partial L}{\partial \theta}$$

- Contudo, quando implementamos a descida de gradiente usamos a versão discreta

$$\theta_{t+1} = \theta_t - \eta \frac{\partial L[\theta_t]}{\partial \theta}$$



Regularização Implícita

- Foi demonstrado que a versão discreta da descida de gradiente aproxima a versão contínua com a seguinte função de custo

$$L_{GD}(\theta) = L(\theta) + \frac{\eta}{4} \left\| \frac{\partial L}{\partial \theta} \right\|^2$$

De modo que a solução é repelida de lugar nos quais a segunda norma do gradiente é grande (alta declividade)

Regularização Implícita

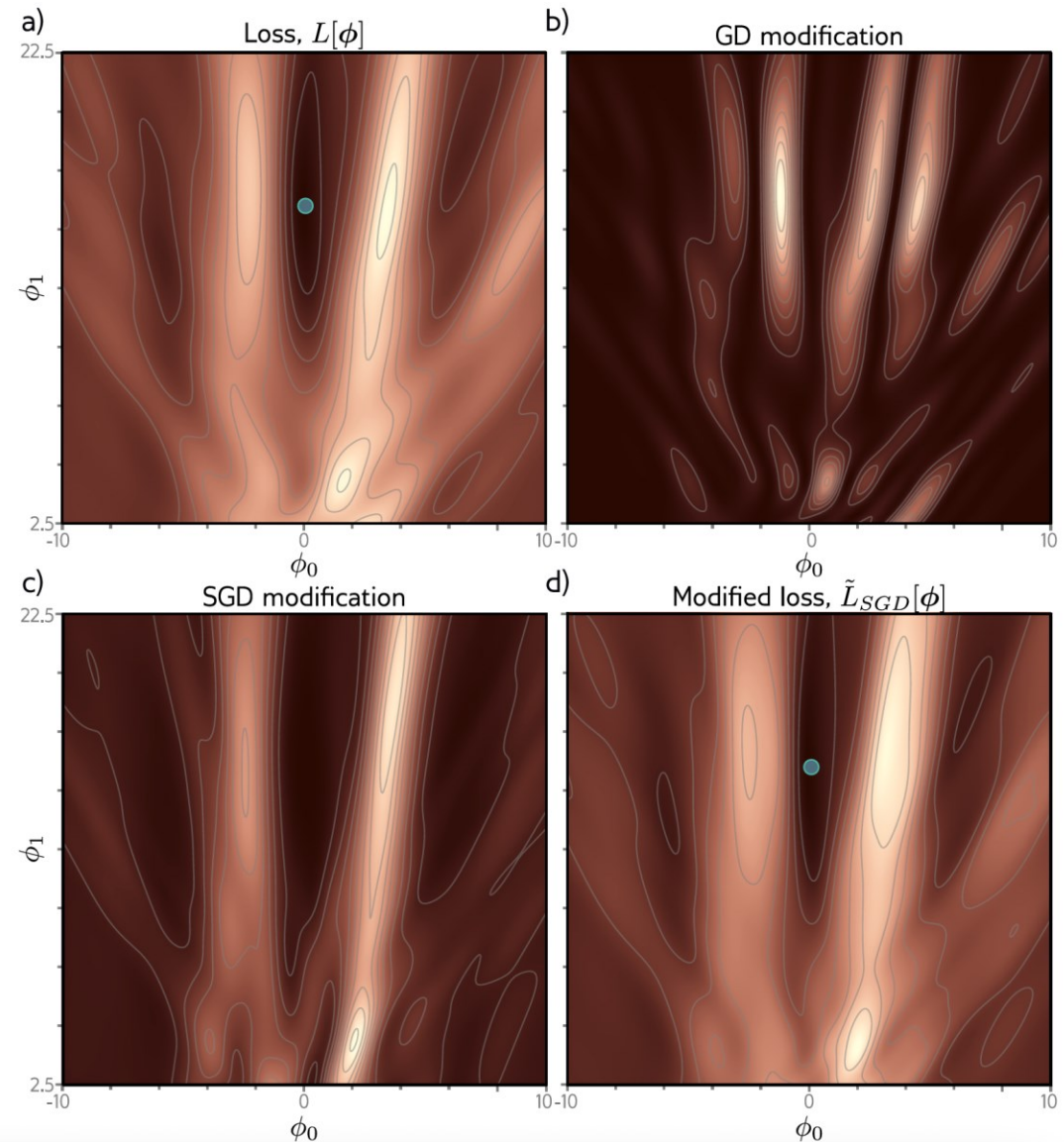
- Foi demonstrado que a versão discreta da descida de gradiente estocástica aproxima a versão contínua com a seguinte função de custo

$$L_{SGD}(\theta) = L(\theta) + \frac{\eta}{4} \left\| \frac{\partial L}{\partial \theta} \right\|^2 + \frac{\eta}{4B} \sum_{b=1}^B \left\| \frac{\partial L_b}{\partial \theta} - \frac{\partial L}{\partial \theta} \right\|^2$$

- De modo que a solução é repelida de lugar nos quais a segunda norma do gradiente é grande (alta declividade)
- O segundo termo de regularização corresponde a variância do gradiente nos diferentes batchs
 - Em outras palavras, SGD favorece regiões cujo gradiente é estável

Regularização Implícita

- Ambas modificações não mudam as posições dos mínimos
 - Mas podem mudar as trajetórias e qual mínimo local encontrado



Referências:

- Sugere-se ***fortemente*** a leitura de:
 - Capítulo 9 de Understanding Deep Learning
 - <https://udlbook.github.io/udlbook/>