

Regularization

$$\begin{aligned} & \arg\min_{\theta} L(D; \theta) + \underbrace{\lambda R(\theta)}_{\text{Regularizer}} \\ & \equiv \arg\max_{\theta} P(\theta | D) = P(D | \theta) \underbrace{P(\theta)}_{\exp(-\lambda R(\theta))} \end{aligned}$$

L-2 regularization

$$\arg\min_{\underline{w}} \underbrace{\|\underline{y} - X\underline{w}\|_2^2}_{\text{Loss function}} + \lambda \underbrace{\|\underline{w}\|_2^2}_{\text{L-2 norm of parameter}}$$

$$\arg \min_{\theta} \underbrace{\sum_{i=1}^n \|y - f(x; \theta)\|}_{\hat{y}} + \lambda \|\theta\|_2^2$$

$$\rightarrow \theta_{t+1} = \theta_t - \alpha \left(\frac{\partial L}{\partial \theta} + 2\lambda \theta_t \right)$$

Neural network call L-2 regularization

as weight decay $\cdot (Wx + b)$

weight decay parameter $= \lambda$

$$\|\theta\|^2 = \theta^T \theta$$

$$= \theta^T I_{n \times n} \theta$$

$$\frac{\partial}{\partial \theta} \theta^T I_{n \times n} \theta = 2\theta$$

Parameters $\rightarrow \theta$

trained by gradient descent or SGD

Hyper parameters

the Loss function should be fixed

$\alpha, \lambda, \text{Batch size}$

(3) Contrastive divergence

① Grid search

$$\lambda = [0.0, 0.01, 0.001, 0.1, 0.2, \dots, 0.9]$$

a) Train

b) find validation accuracy

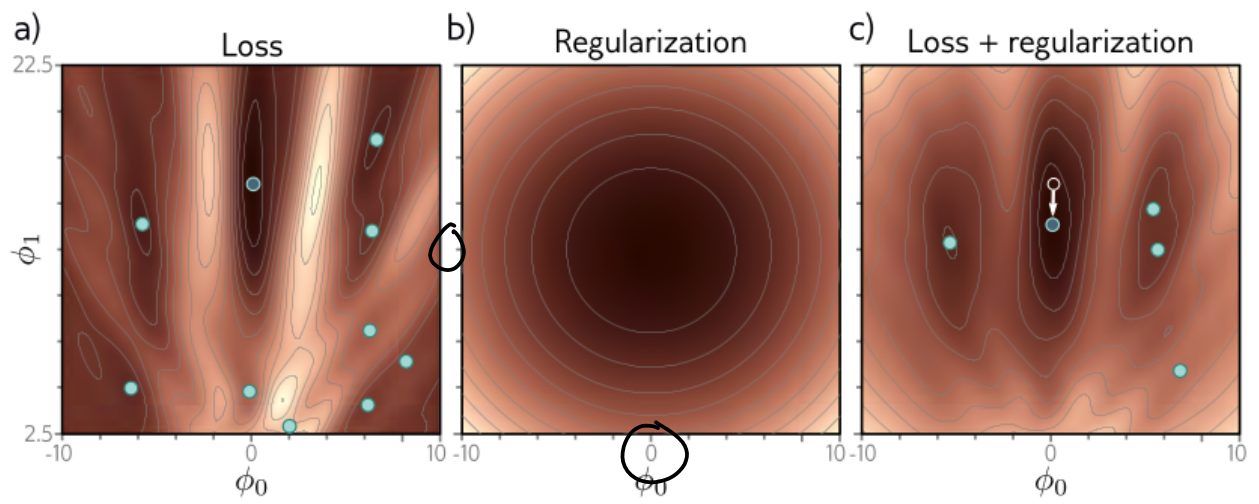
c) Pick the best set of hyper parameters

② Random search

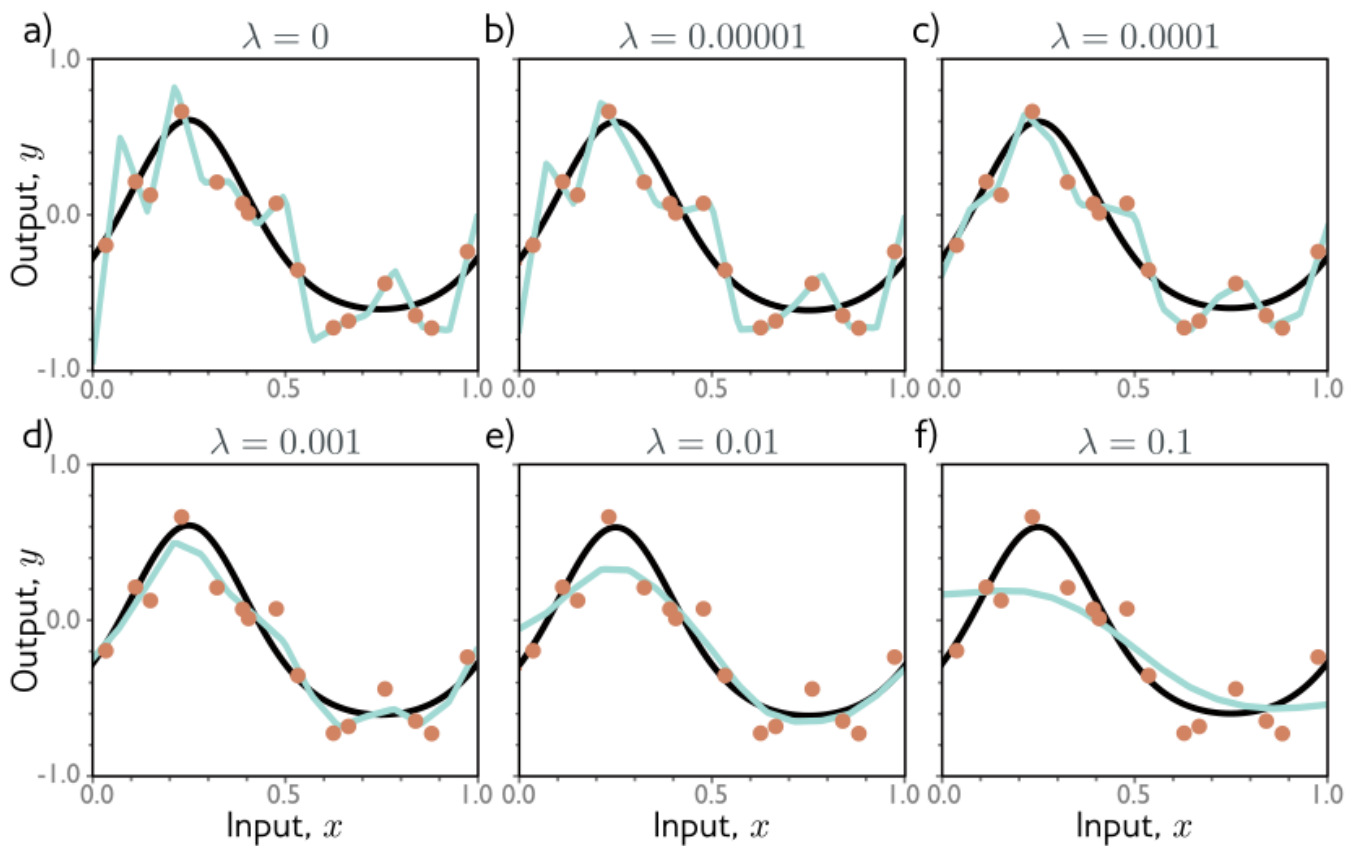
$$\lambda \sim U[0.0001, 0.99]$$

a) Train b) find validation c) Pick the best

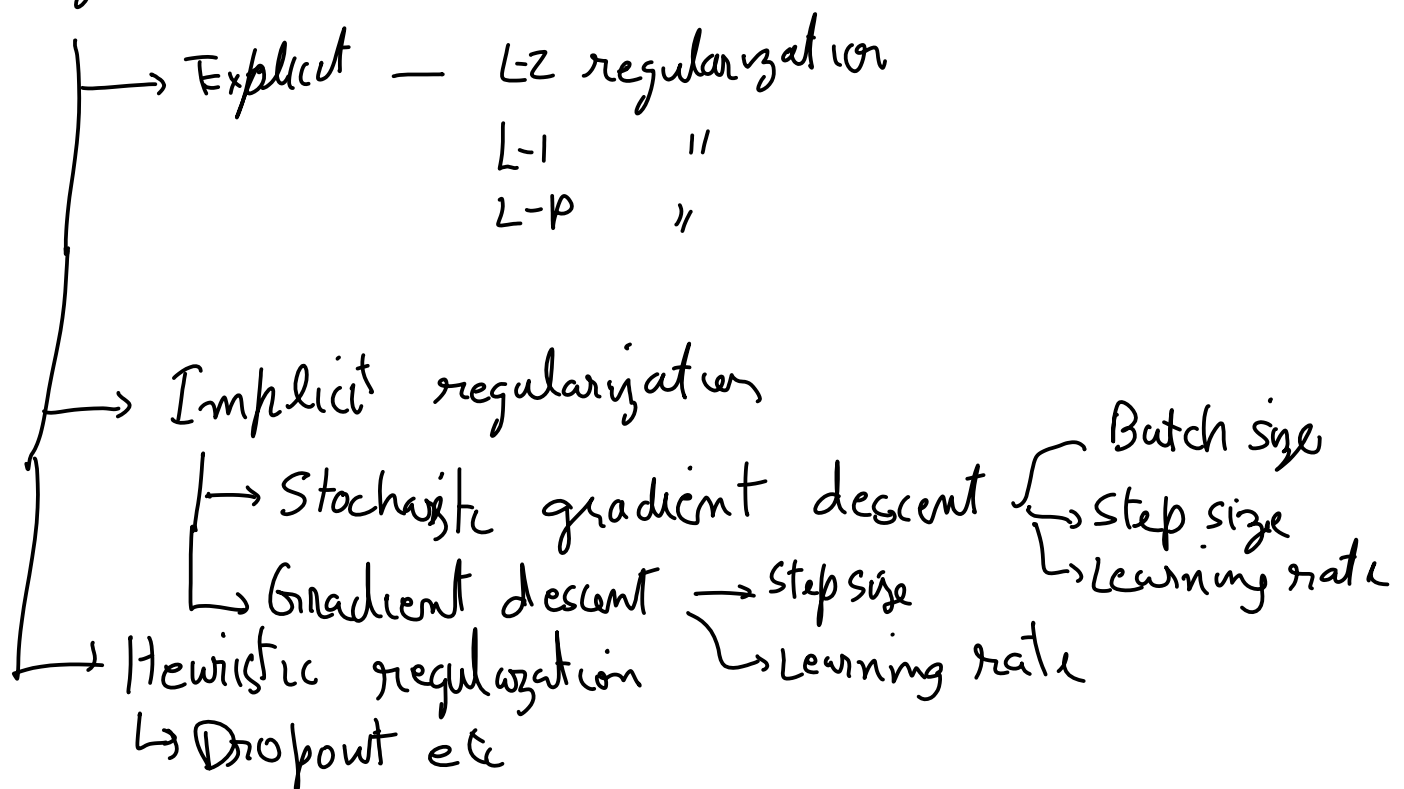
Explicit regularization



L2 Regularization



Regularization



GD

Gradient descent induce regularization

$$\theta_{t+1} = \theta_t - \alpha \left. \frac{\partial L}{\partial \theta} \right|_{\theta = \theta_t}$$

Ideally GD should converge to Local Minima

Gradient flow

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

$$\dot{\theta} = -g(\theta(t))$$

Ideal

Actually: $\theta_{t+1} = \theta_t - g(\theta(t))$

Corrected
Gradient flow

$$\dot{\theta} = -\tilde{g}(\theta(t)) = -g(\theta(t)) - \epsilon \underbrace{g_1(\theta(t))}_{\text{error term}}$$

$$-\frac{\partial \tilde{L}_{GD}}{\partial \theta} = \frac{\partial L}{\partial \theta} + (\epsilon \text{ term})$$

regularizer



$$\theta_{t+1} = \theta_t - \alpha \frac{\partial L}{\partial \theta}$$

By Taylor series expansion

$$\underline{\theta(t+\epsilon)} = \underline{\theta(t)} + \epsilon \frac{\partial \underline{\theta(t)}}{\partial t} + \frac{\epsilon^2}{2} \frac{\partial^2 \underline{\theta(t)}}{\partial t^2} + \dots$$

$$\frac{\partial \underline{\theta}}{\partial t}(t) = -\tilde{g}(\theta(t)) = -g(\theta(t)) - \epsilon g_1(\theta(t)) \quad \text{--- (1)}$$

$$\frac{\partial^2 \underline{\theta}}{\partial t^2}(t) = -\frac{\partial g(\theta(t))}{\partial \theta} \frac{\partial \underline{\theta}}{\partial t} - \epsilon \left[\frac{\partial}{\partial t} g_1(\theta(t)) \right] \quad \text{--- (2)}$$

$$\theta(t+\epsilon) = \theta(t) - \epsilon g(\theta(t)) - \epsilon^2 g_1(\theta(t))$$

$$- \underbrace{\frac{\epsilon^2}{2} \left[\frac{\partial g(\theta(t))}{\partial \theta} \frac{\partial \underline{\theta}}{\partial t} \right]}_{\text{(2)}} - \underbrace{\left[\frac{\epsilon^3}{2} \frac{\partial}{\partial t} g_1(\theta(t)) \right]}_{O(\epsilon^3)}$$

If ϵ is very small we can ignore terms with ϵ^3

To make the equation as close as possible to

$$\theta(t+\epsilon) = \theta(t) - \epsilon g(\theta(t))$$

we need $g_1(\theta(t)) = -\frac{1}{2} \frac{\partial g(\theta(t))}{\partial \theta} \frac{\partial \theta}{\partial t}$

Corrected gradient flow

$$\dot{\theta} = -g(\theta) - \epsilon g_1(\theta)$$

$$= -g(\theta) + \frac{\epsilon}{2} \frac{\partial}{\partial \theta} g(\theta) \frac{\partial \theta}{\partial t}$$

$$= -\frac{\partial L}{\partial \theta} + \frac{\epsilon}{2} \frac{\partial^2 L(\theta)}{\partial \theta^2} \frac{\partial \theta}{\partial t} \rightarrow \boxed{\frac{\partial \theta}{\partial t} = -g(\theta)} = -\frac{\partial L}{\partial \theta}$$

$$\dot{\theta} = -\frac{\partial L}{\partial \theta} - \frac{\epsilon}{2} \left[\frac{\partial^2 L(\theta)}{\partial \theta^2} \right] \frac{\partial L}{\partial \theta}$$

← scalar

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

$$L_{GD} = L + \underbrace{\frac{\alpha}{4} \left\| \frac{\partial L}{\partial \theta} \right\|_2^2}_{\text{Regularizer}}$$

Gradient descent
implicitly introduces
a regularizer of the size
proportional to step size α
and squared norm of gradients $\left\| \frac{\partial L}{\partial \theta} \right\|_2^2$

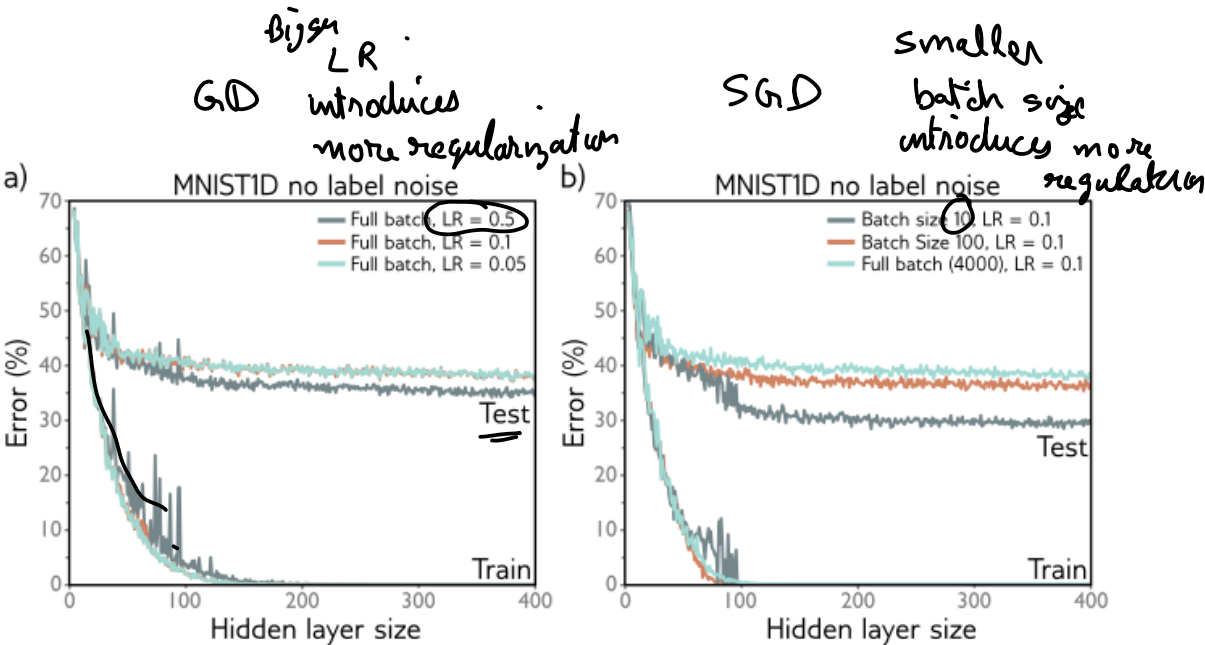
$$\begin{aligned} \frac{\partial L}{\partial \theta} &= g \\ \frac{\partial}{\partial \theta} \left(\|g\|^2 \right) &= 2g \frac{\partial g}{\partial \theta} \\ &= 2 \frac{\partial L}{\partial \theta} \frac{\partial^2 L}{\partial \theta^2} \end{aligned}$$

SGD introduces additional regularizer
to GD

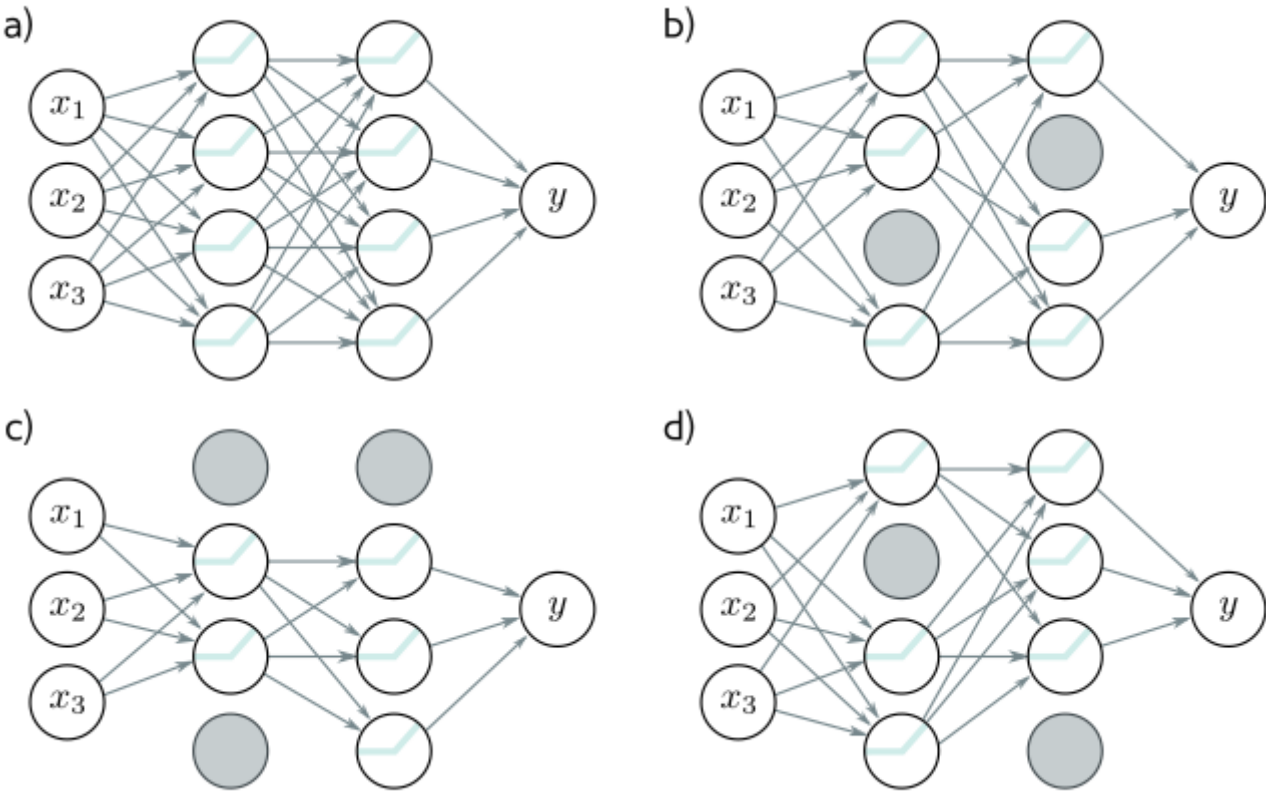
that is proportional to the variance of batch
gradients.

$$L_{\text{SGD}} = L_{\text{GD}} + \epsilon^2 \sum_{i=1}^m \left\| \left(\frac{\partial L_B}{\partial \theta} \right) - \frac{\partial L}{\partial \theta} \right\|_2^2$$

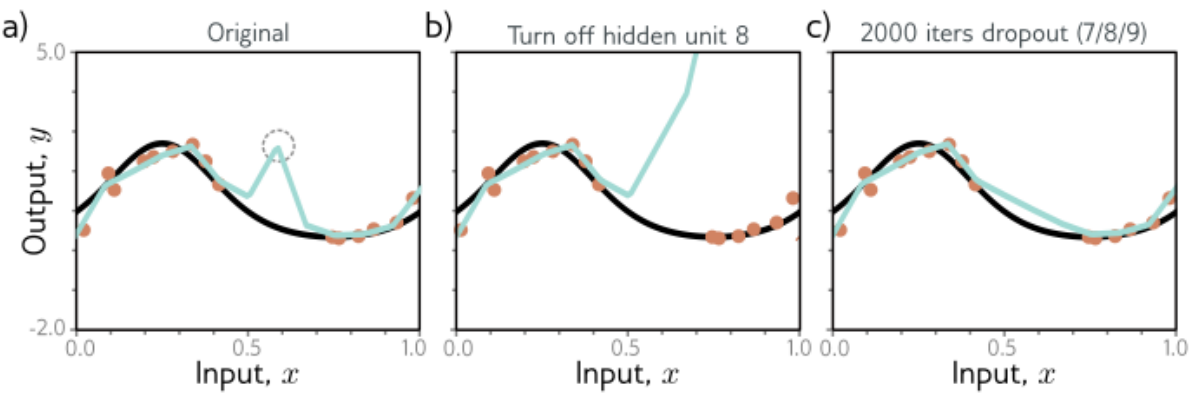
Effect of batch size and learning rate



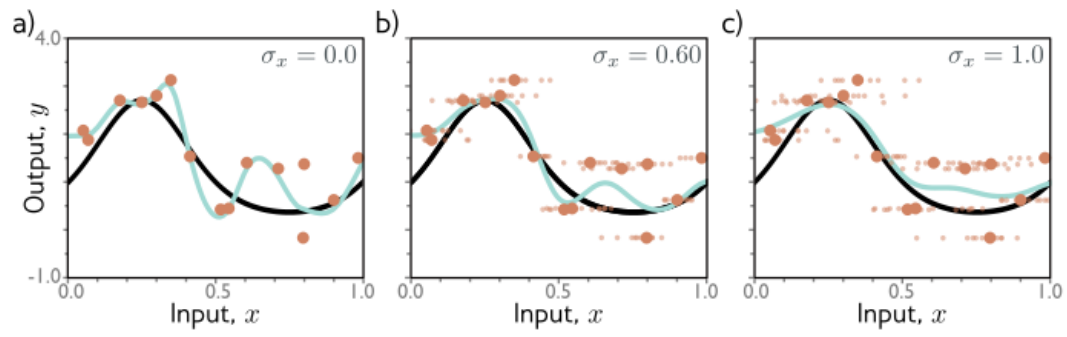
Dropout



Effect of dropout



Adding noise to each batch



Data augmentation

