# Training the Network

layer 3
feature maps

$G_n = f(L_n; \omega_1, \dots \omega_K)$ $\longrightarrow$ $\ell(G_n; W) = \ell_n$

$\omega_1\ \omega_2 \qquad\qquad \omega_K$

layer 3 filters

layer 2
feature maps

$L_n = f(M_n; \Psi_1, \dots \Psi_K)$

$\Psi_1\ \Psi_2 \qquad\qquad \Psi_K$

layer 2 filters

input image

layer 1
feature maps

$M_n = f(I_n; \phi_1, \dots \phi_K)$

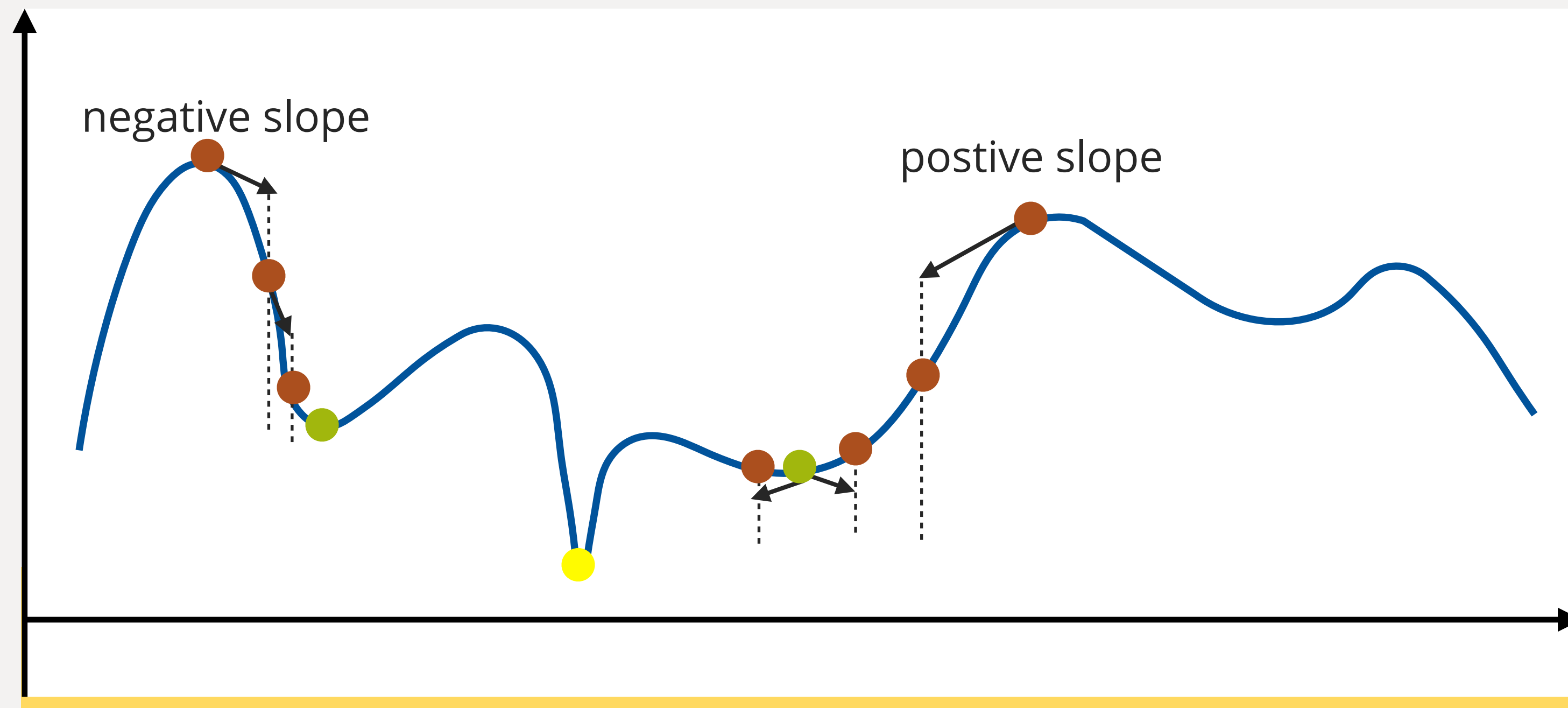$\phi_1\ \phi_2 \qquad\qquad \phi_K$

layer 1 filters

$I_n$

## How the Model Learns

Assume we have labeled images
$\{I_n, y_n\}\ n = 1, N$
$I_n$ is image $n$
$y_n \in \{+1, -1\}$ is associated label

Risk function of model parameters

$$E(\Phi, \Psi, \Omega, W) = 1/N \sum_{n=1}^{N} loss(y_n, \ell_n)$$

Find model parameters $\hat{\Phi}, \hat{\Psi}, \hat{\Omega}, \hat{W}$
that minimize $E(\Phi, \Psi, \Omega, W)$

# Gradient Descent

$$\Theta = \{ \Phi, \Psi, \Omega, W \}$$



negative slope

postive slope

$$\Theta_{t+1} = \Theta_t - \alpha \nabla_\Theta E(\Theta_t)$$

multi-dimensional "slope"

# Gradient Descent $\longrightarrow$ Stochastic Gradient Descent

$$\Theta_{t+1} = \Theta_t - \alpha \nabla_\Theta E(\Theta_t) \qquad \Theta_{t+1} = \Theta_t - \alpha \nabla_\Theta \hat{E}(\Theta_t)$$

$$\hat{E}_t(\Phi, \Psi, \Omega, W) = 1/|S_t| \sum_{n \in S_t} loss(y_n, \ell_n)$$

random subset
of data

## Massive $N$

- Choose a **random** data subset
- Estimate gradient by data point
- Update parameters using gradient from random subset
- Leads to similar solutions at faster rate