

CS 273A: Machine Learning

Winter 2021

Lecture 15: Latent-Space Models

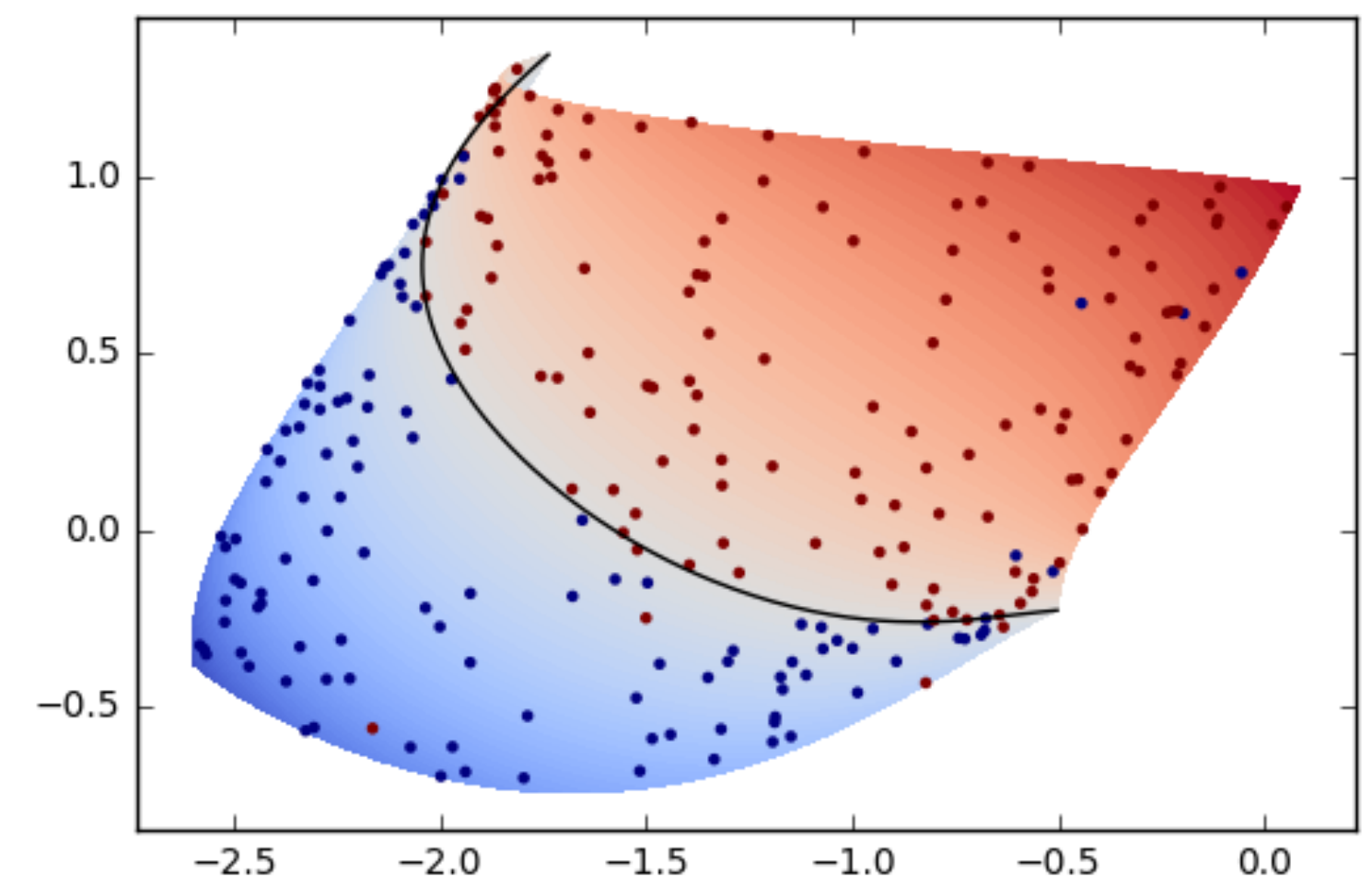
Roy Fox

Department of Computer Science

Bren School of Information and Computer Sciences

University of California, Irvine

All slides in this course adapted from Alex Ihler & Sameer Singh



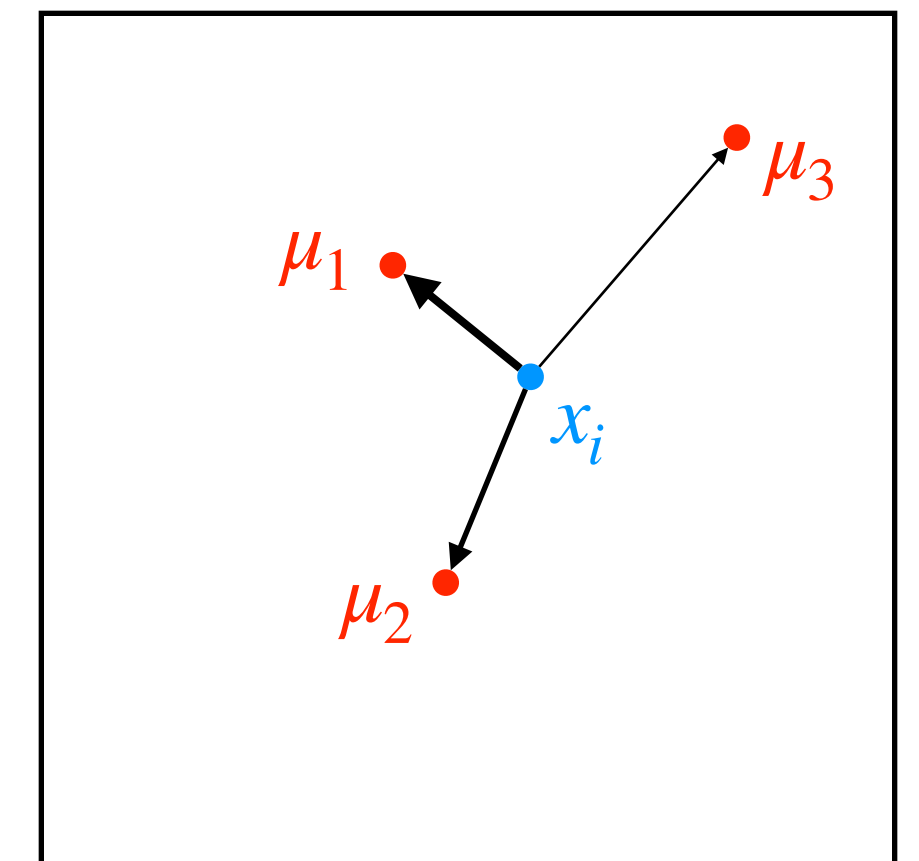
Today's lecture

Gaussian Mixture Models

Dimensionality Reduction

Mixture Models

- k -Means assigns each instance to **one cluster**
 - ▶ Could it be assigned to **another cluster** equally well? Almost equally?
 - ▶ **Hard assignment** $f : x \mapsto c$ loses information on:
 - Which clusters are “**close seconds**”
 - **Uncertainty** = how sure are we of the assignment
- **Mixture Model** = **prior** over clusters $p(c)$ + **distribution** in each cluster $p(x | c)$
 - ▶ \implies **Posterior** $p(c | x)$ = probabilistic (**soft**) assignment of x to c



Gaussian Mixture Models (GMMs)

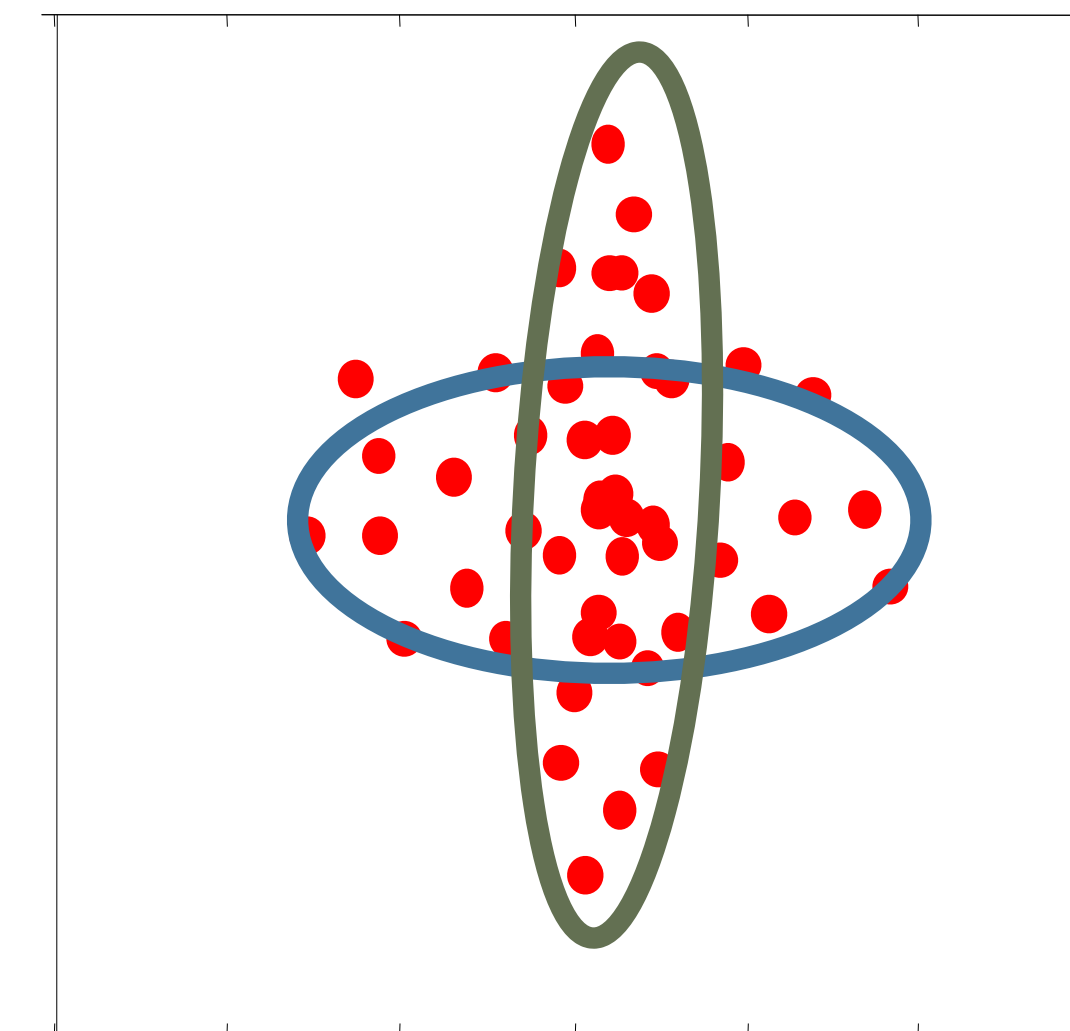
- Each cluster is modeled by a **Gaussian** $p(x | c) = \mathcal{N}(x; \mu_c, \Sigma_c)$
 - Σ_c allows **non-isotropic** clusters \implies **weighted** Euclidean distance
- **Mixture** = distribution over Gaussians is given by a probability vector $p(c)$
- **Generative model** = we can sample $p(x)$:

- Sample $z \sim p(c)$

- Sample $x \sim p(x | c = z)$

we don't output z , it is "latent" = hidden
 \implies **can be any of them**

- Probability of this x : $\sum_c p(c = z)p(x | c = z) = \sum_c p(c, x) = p(x)$



Multivariate Gaussian distributions

$$\mathcal{N}(x; \mu, \Sigma) = (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$

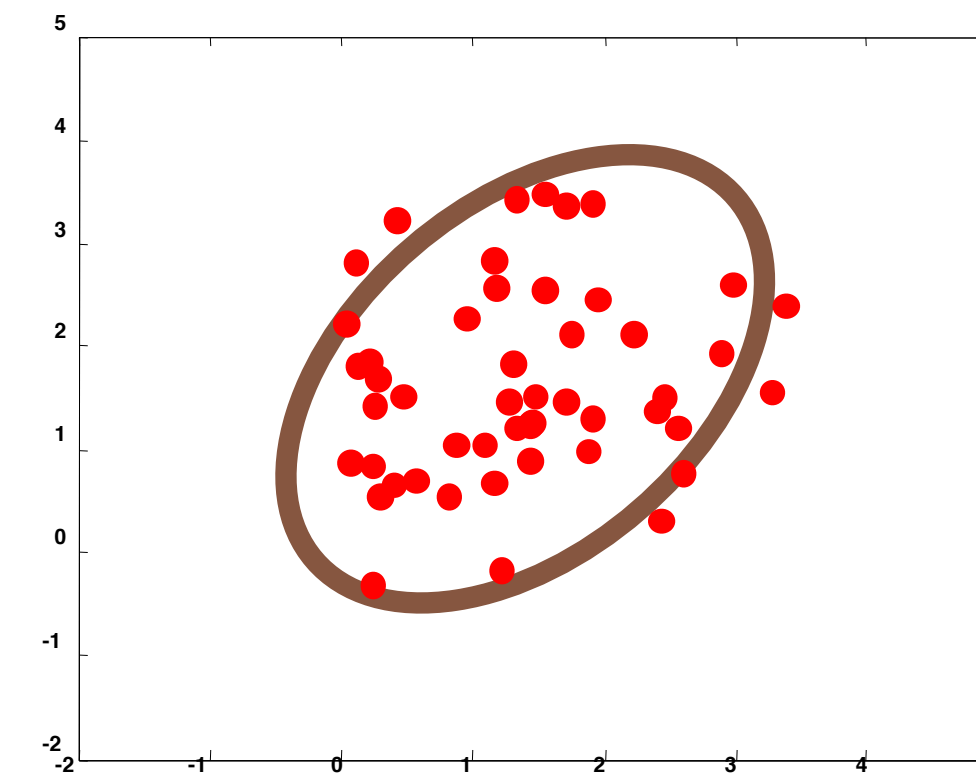
- For data points $\{x_i\}$, maximum log-likelihood estimator of μ, Σ :

$$\nabla_{\mu} \sum_i \log \mathcal{N}(x_i; \mu, \Sigma) = \frac{1}{2} \sum_i (x_i - \mu)^\top \Sigma^{-1} = 0$$

$$\implies \mu = \frac{1}{m} \sum_i x_i$$

$$\nabla_{\Sigma^{-1}} \sum_i \log \mathcal{N}(x_i; \mu, \Sigma) = -\frac{1}{2} \sum_i \left((x_i - \mu)(x_i - \mu)^\top - \Sigma \right) = 0$$

$$\implies \Sigma = \frac{1}{m} \sum_i (x_i - \mu)(x_i - \mu)^\top$$



matrix calculus identity:
 $\nabla_{\Sigma^{-1}} \log |\Sigma|^{-1} = \Sigma$

Training GMMs

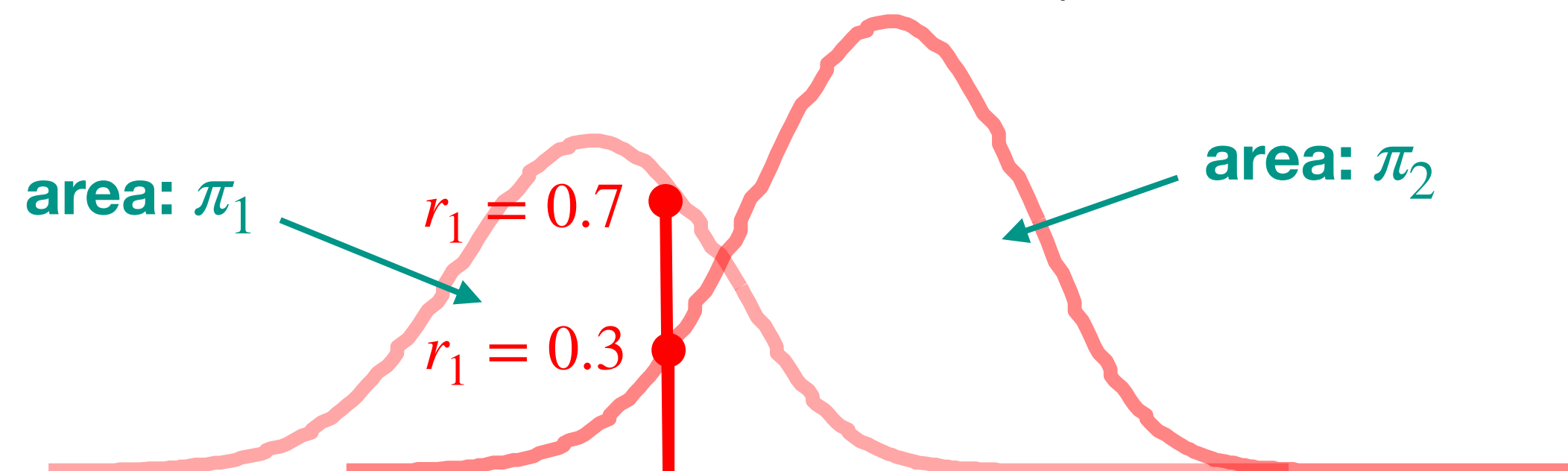
- k -Means:
 - ▶ Assign data points to clusters z_i
 - ▶ Update each cluster's parameters μ_c
- A “soft” version of k -Means: Expectation–Maximization (EM) algorithm
 - ▶ Find a “soft” assignment $p(c | x)$
 - ▶ Update model parameters $p(c), p(x | c)$
- The EM algorithm is extremely general, GMMs are a very special case

Expectation–Maximization: E-step

- **Initialize** model parameters $\pi_c = p(c), \mu_c, \Sigma_c$
- **E-step (Expectation)**: [why “expectation”? comes from the general EM algorithm]
 - For each data point x_i , use **Bayes' rule** to compute:

$$r_{ic} = p(c | x_i) = \frac{p(c)p(x_i | c)}{\sum_{\bar{c}} p(\bar{c})p(x_i | \bar{c})} = \frac{\pi_c \mathcal{N}(x_i; \mu_c, \Sigma_c)}{\sum_{\bar{c}} \pi_{\bar{c}} \mathcal{N}(x_i; \mu_{\bar{c}}, \Sigma_{\bar{c}})}$$

- High weight to clusters that are **likely a-priori**, or in which x_i is **relatively probable**



Expectation–Maximization: M-step

- Given assignment probabilities r_{ic}
- M-step (Maximization):
 - For each cluster c , fit the best Gaussian to the weighted assignment

total weight assigned to cluster c

$$m_c = \sum_i r_{ic}$$

what is $\sum_c m_c$? m

fraction of weight assigned to cluster c

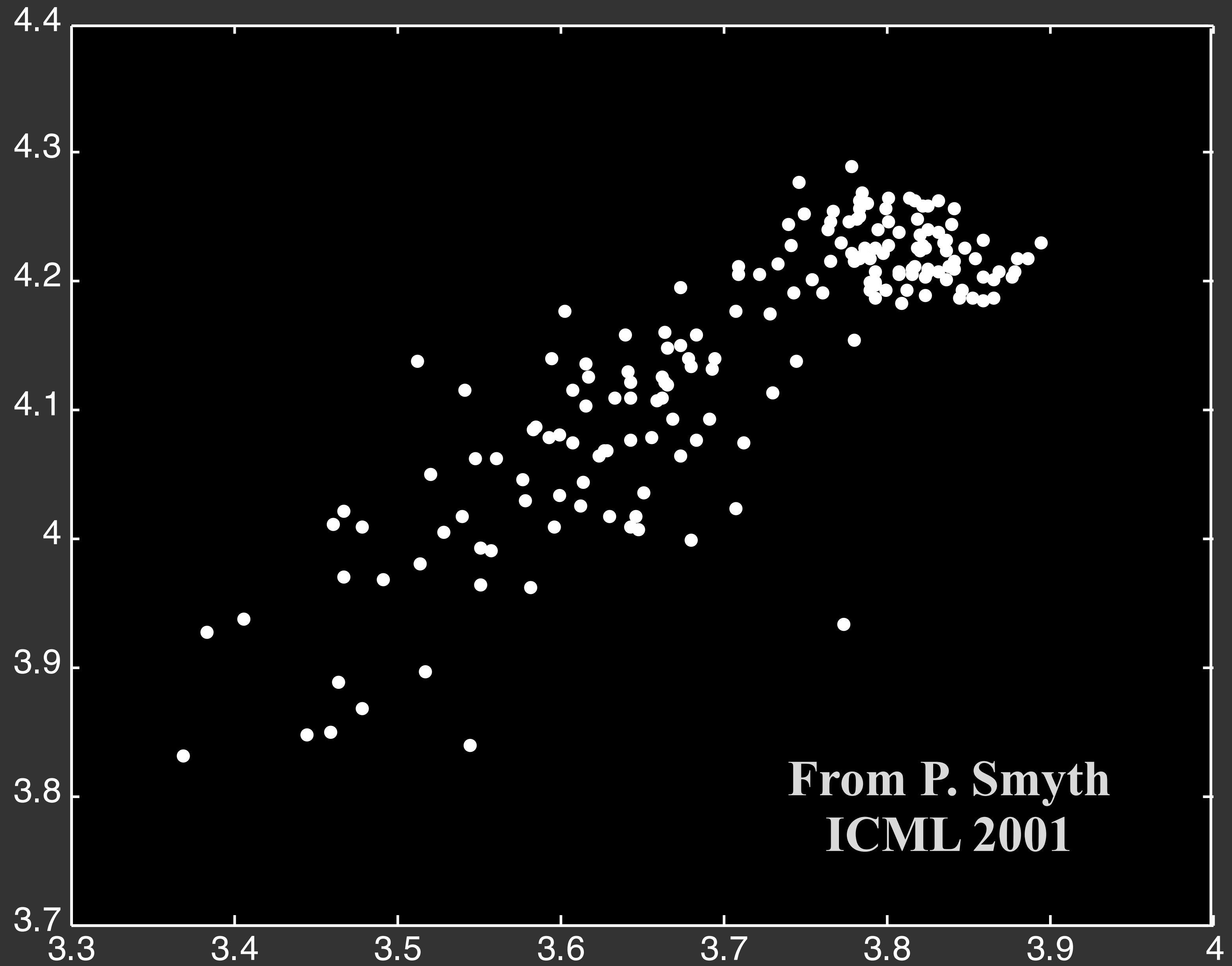
$$\pi_c = \frac{m_c}{m}$$

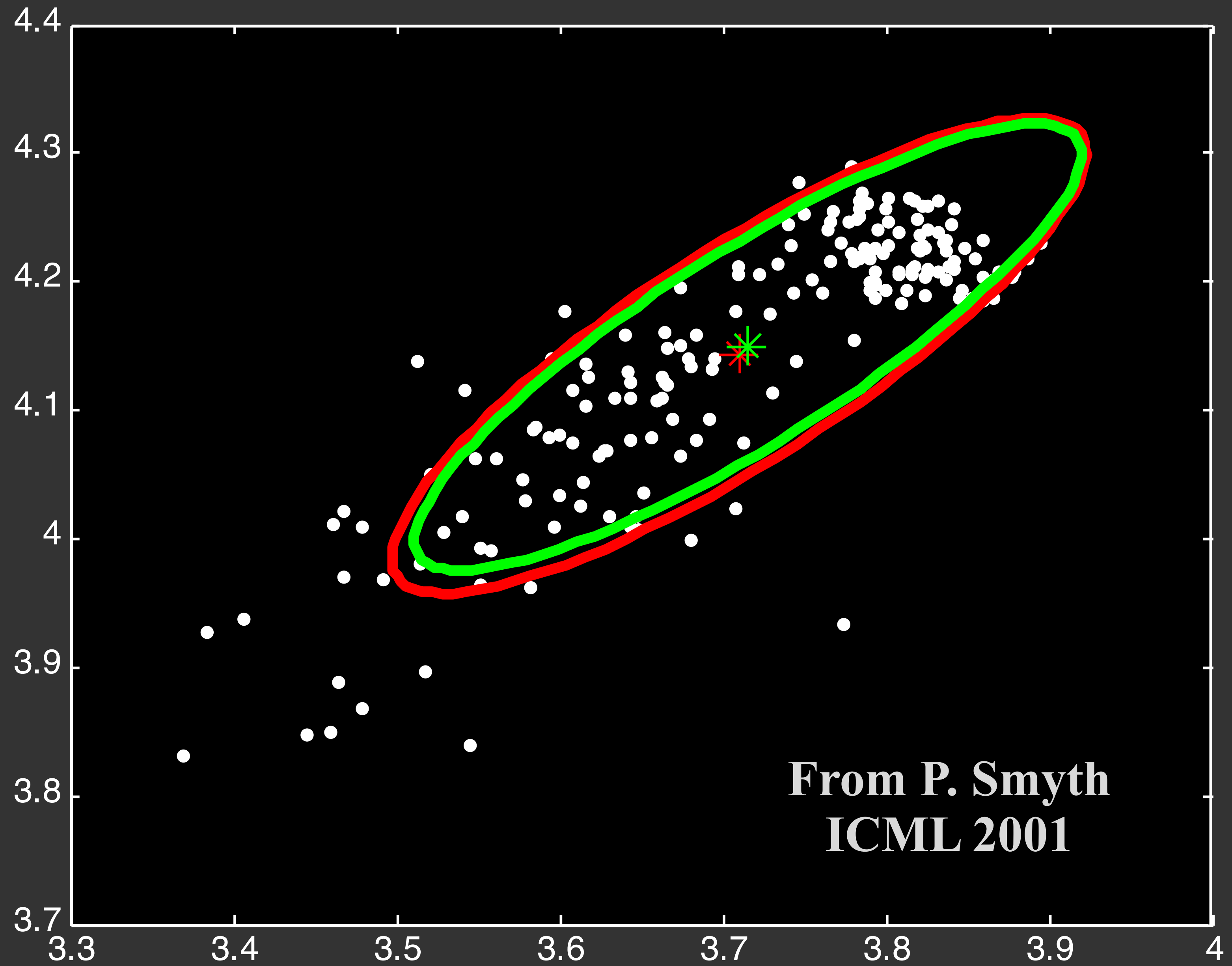
$$\mu_c = \frac{1}{m_c} \sum_i r_{ic} x_i$$

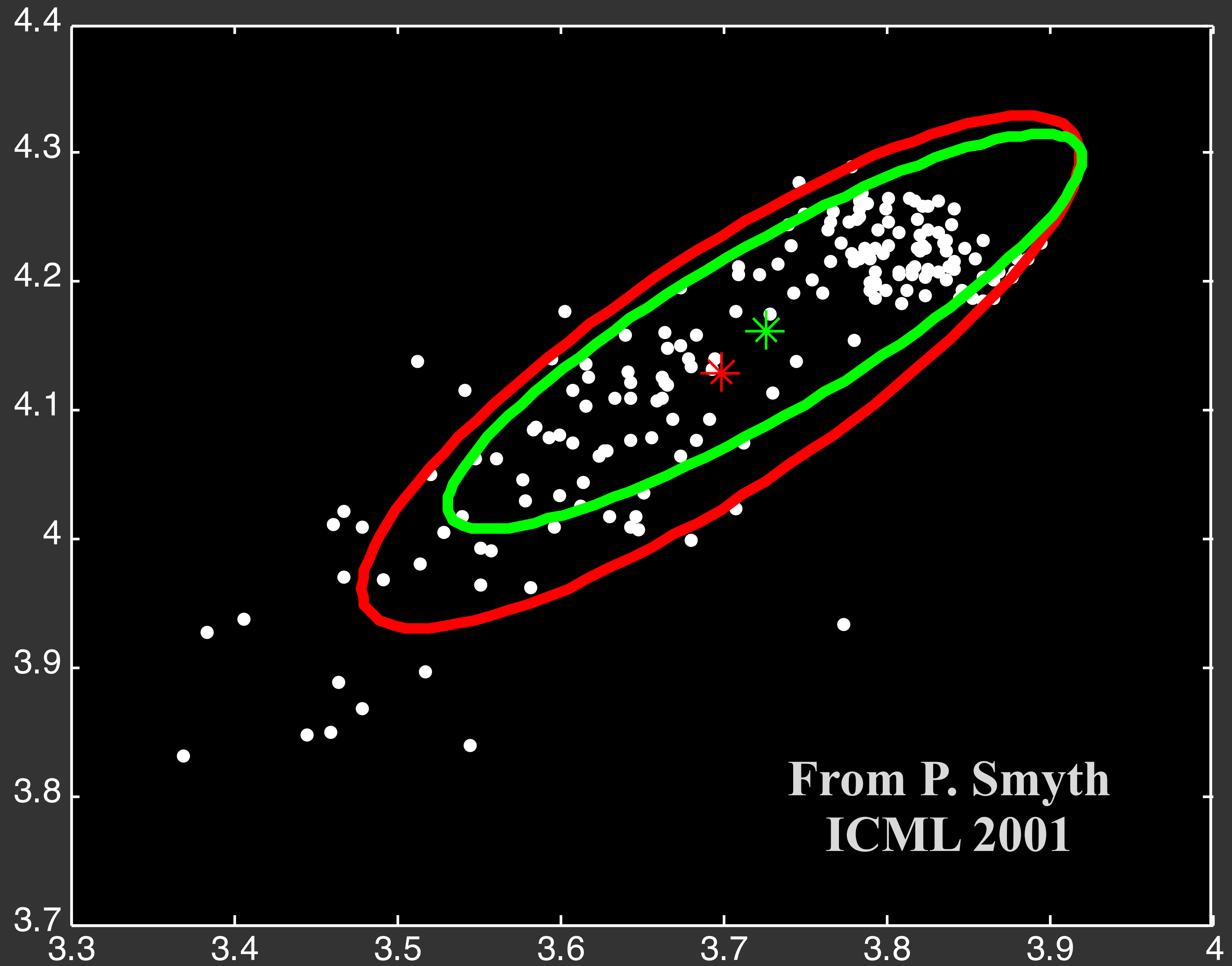
weighted mean of data in cluster c

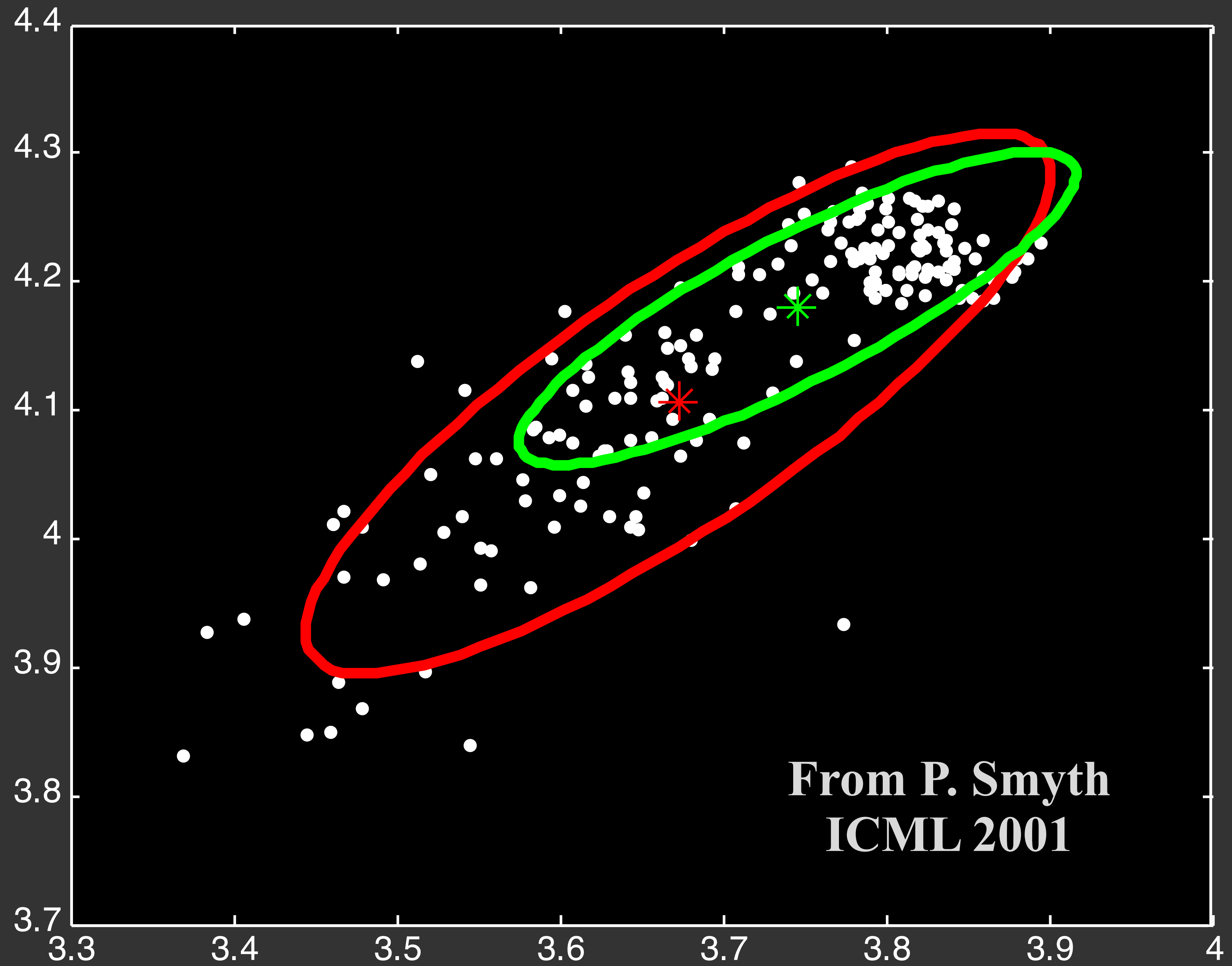
$$\Sigma_c = \frac{1}{m_c} \sum_i r_{ic} (x_i - \mu_c)(x_i - \mu_c)^T$$

weighted covariance of data in cluster c

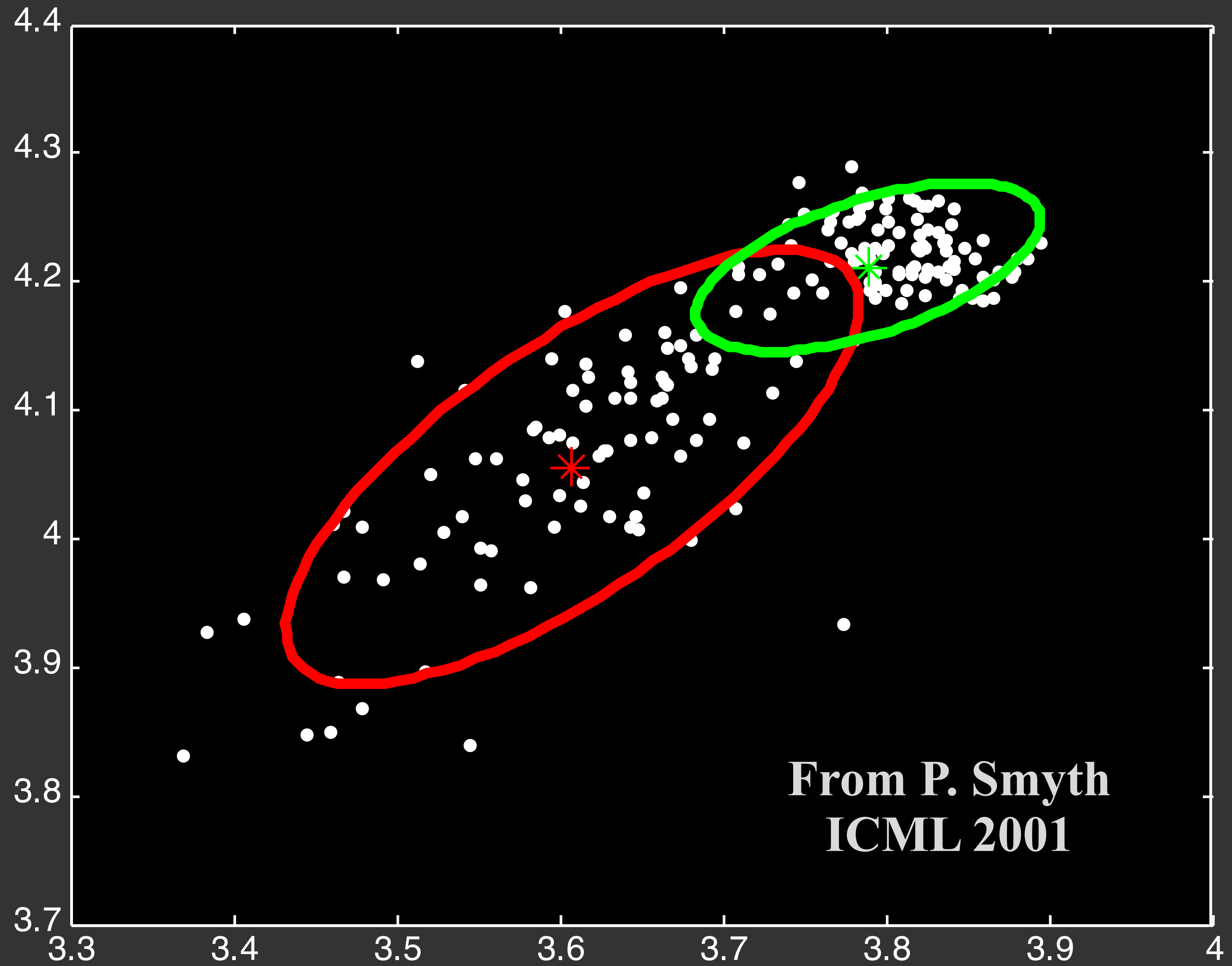




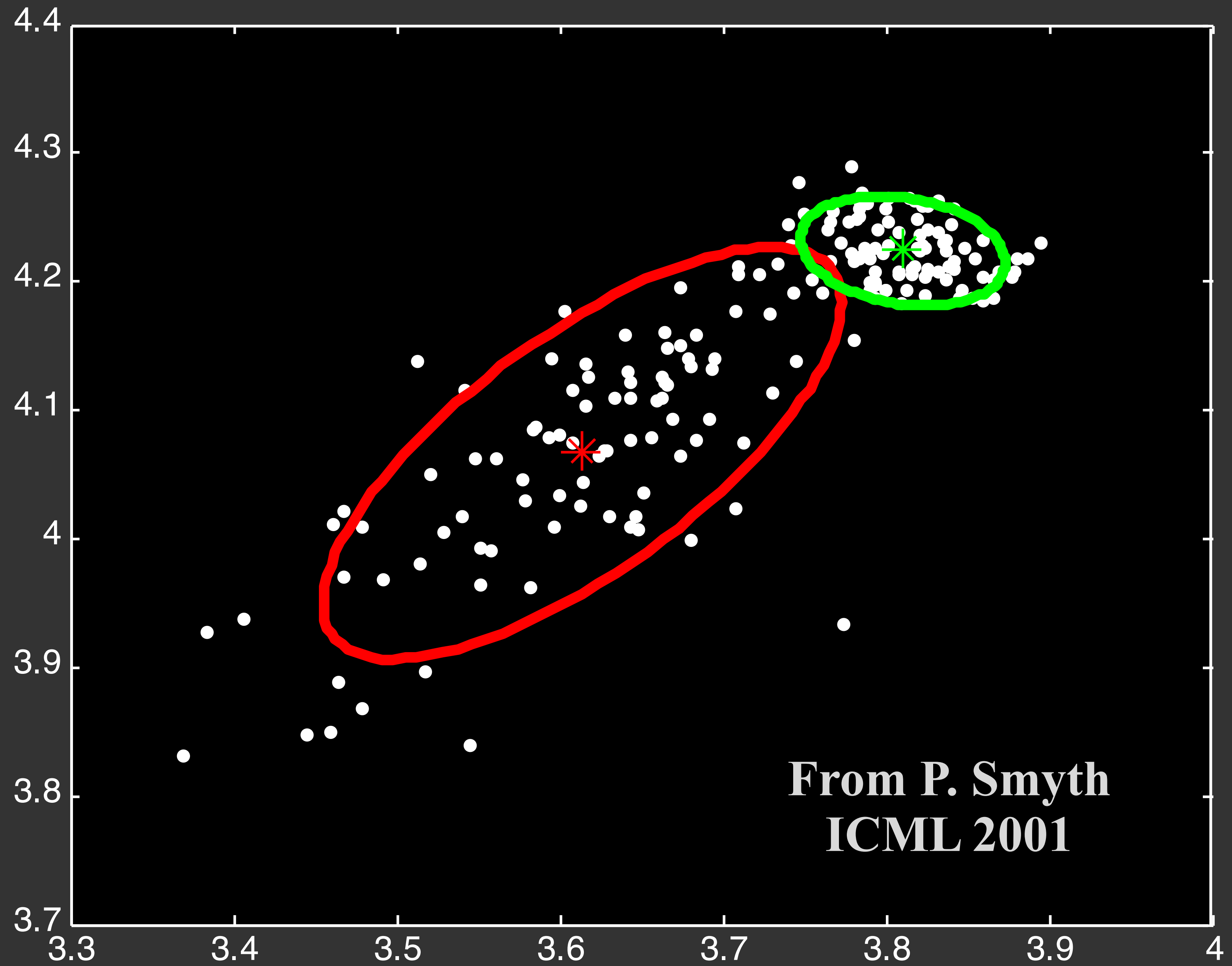




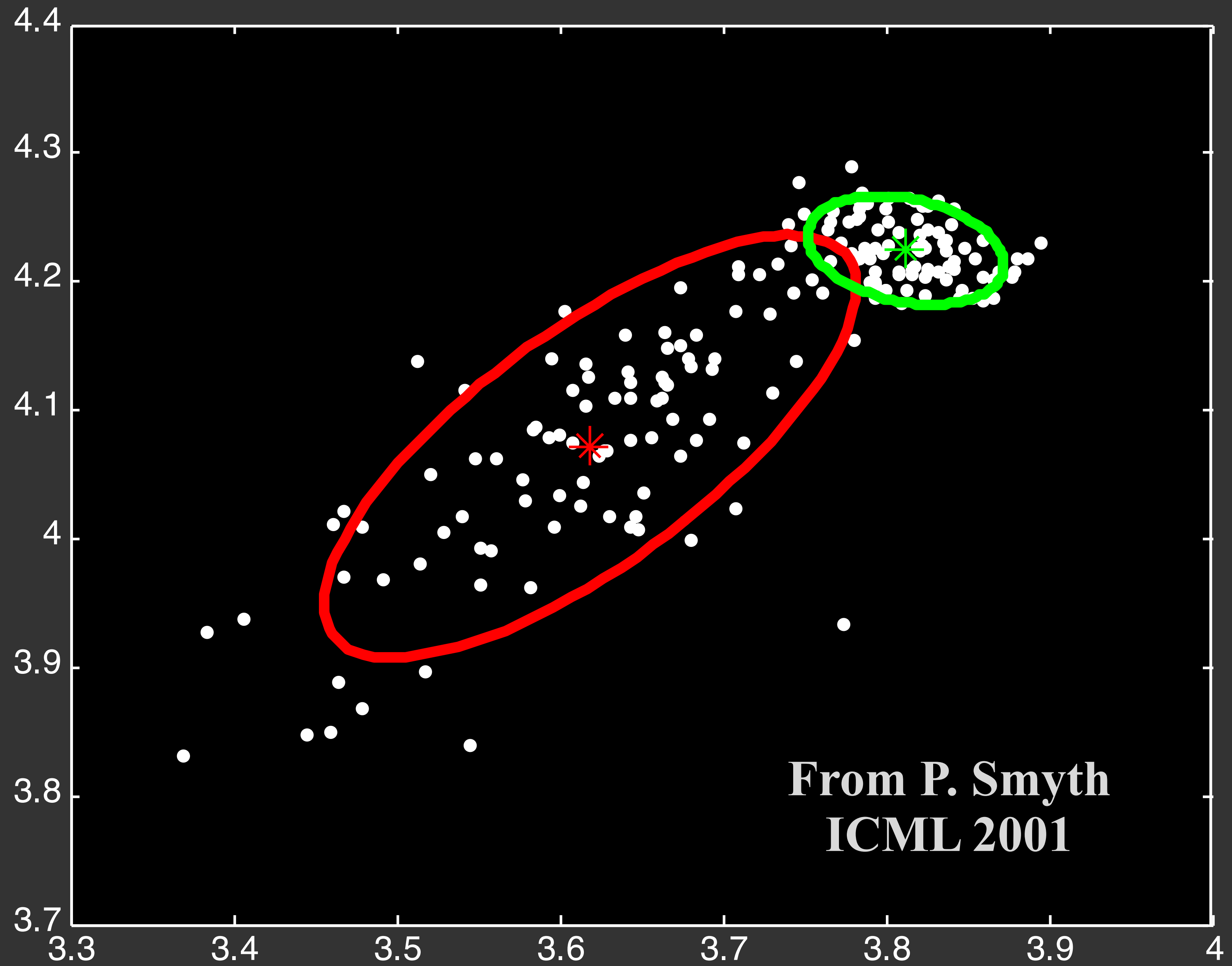
**From P. Smyth
ICML 2001**



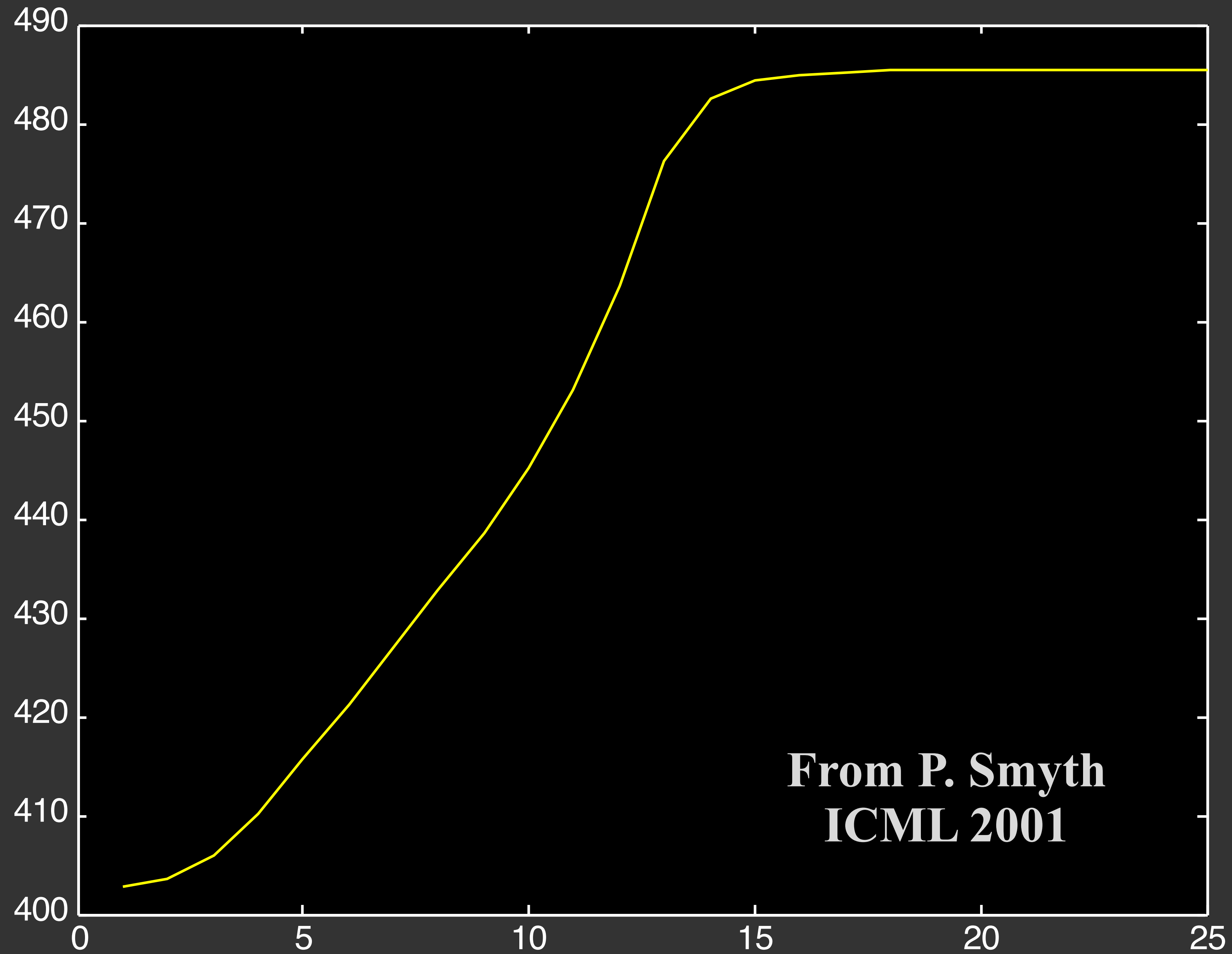
From P. Smyth
ICML 2001



**From P. Smyth
ICML 2001**



From P. Smyth
ICML 2001



**From P. Smyth
ICML 2001**

Demo

- <https://lukapopijac.github.io/gaussian-mixture-model/>

Expectation–Maximization: considerations

- Each iteration of EM is guaranteed to increase the data **log likelihood**

$$\log p(\mathcal{D}) = \sum_i \log p(x_i) = \sum_i \log \sum_c \pi_c \mathcal{N}(x_i; \mu_c, \Sigma_c)$$

**we won't show this
but proof is very insightful!**

- ▶ Convergence **guaranteed** — descends NLL
 - But could be **local optima** \implies **initialization** important
- **Out-of-sample** data: can find **soft assignment** = probabilistic prediction
- Choosing **#clusters**: **regularized** training log-likelihood (as in k -Means)
 - ▶ Or: **validate** log-likelihood on held out data; many clusters \implies **overfitting!**

Recap

- Gaussian Mixture Models (GMMs)
 - Expressive class of generative models $p(x)$
 - Explain variation with latent clusters + cluster distribution
 - Given cluster (= mode), feature values are Gaussian
- Expectation–Maximization (EM)
 - Compute soft assignment probabilities, “responsibility” r_{ic}
 - Update model parameters: mixture π_c , cluster mean and covariance μ_c, Σ_c
 - Ascent on log-likelihood: convergent, but local optima
- Selecting the number of clusters
 - Regularized training log-likelihood, or validation log-likelihood

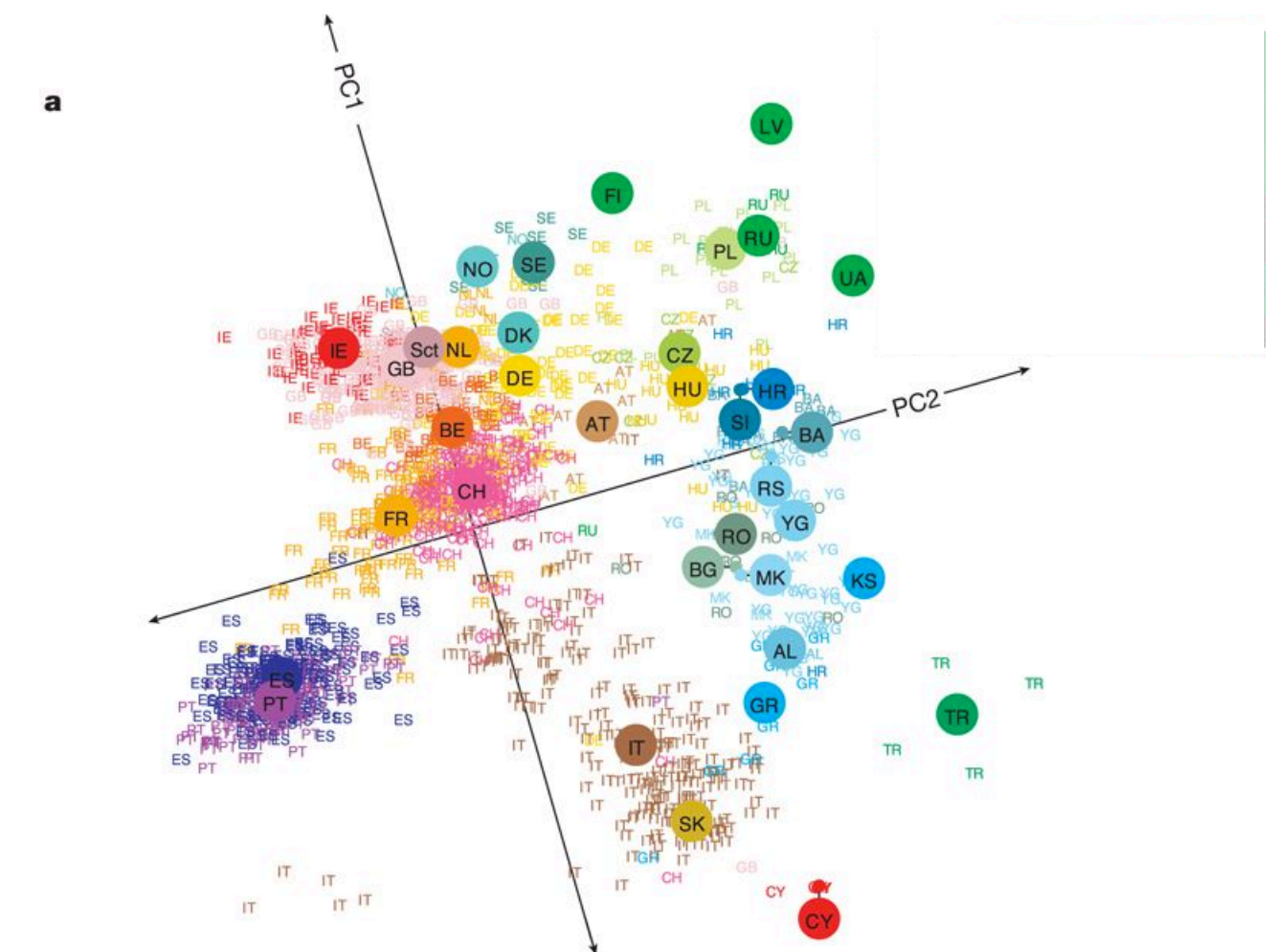
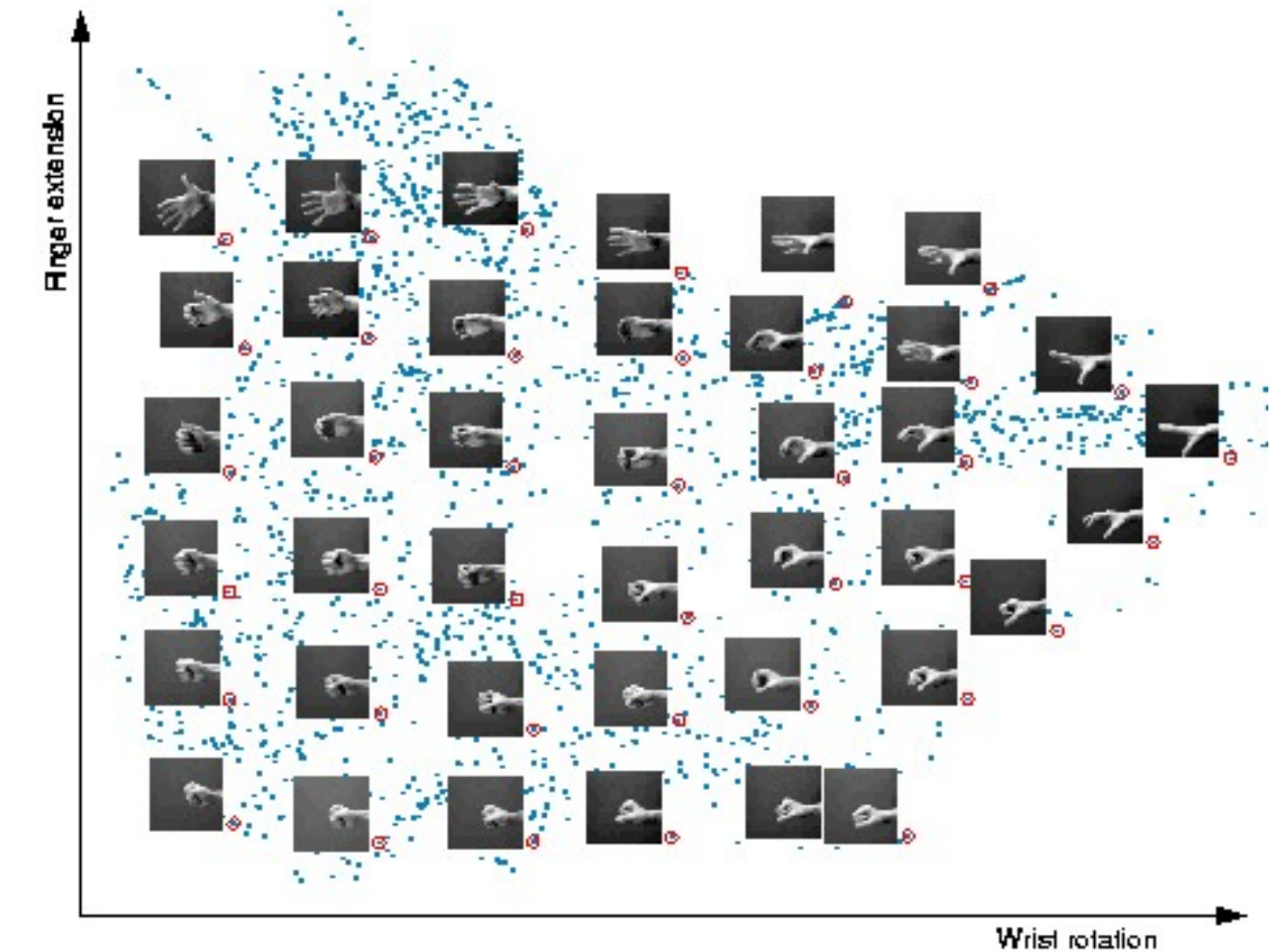
Today's lecture

Gaussian Mixture Models

Dimensionality Reduction

Why reduce dimensionality?

- Data is often **high-dimensional** = many features
 - ▶ **Images** (even at 28x28 pixels)
 - ▶ **Text** (even a “bag of words”)
 - ▶ **Stock prices** (e.g. S&P500)
- Issues with high-dimensionality:
 - ▶ **Computational complexity** of analyzing the data
 - ▶ **Model complexity** (more parameters)
 - ▶ **Sparse data** = cannot cover all combinations of features
 - ▶ Correlated features can be independently **noisy**
 - ▶ Hard to **visualize**



Dimensionality reduction

- With many features, some tend to **change together**
 - Can be **summarized together**
 - Others may have little or **irrelevant change**
- Example: S&P500 → “Tech stocks up 2x, manufacturing up 1.5x, ...”
- **Embed** instances in lower-dimensional space $f: \mathbb{R}^n \mapsto \mathbb{R}^d$
 - Keep dimensions of “interesting” **variability** of data
 - Discard dimensions of **noise** or unimportant variability; or no variability at all
 - Keep “similar” data **close** \implies may preserve **cluster structure**, other insights

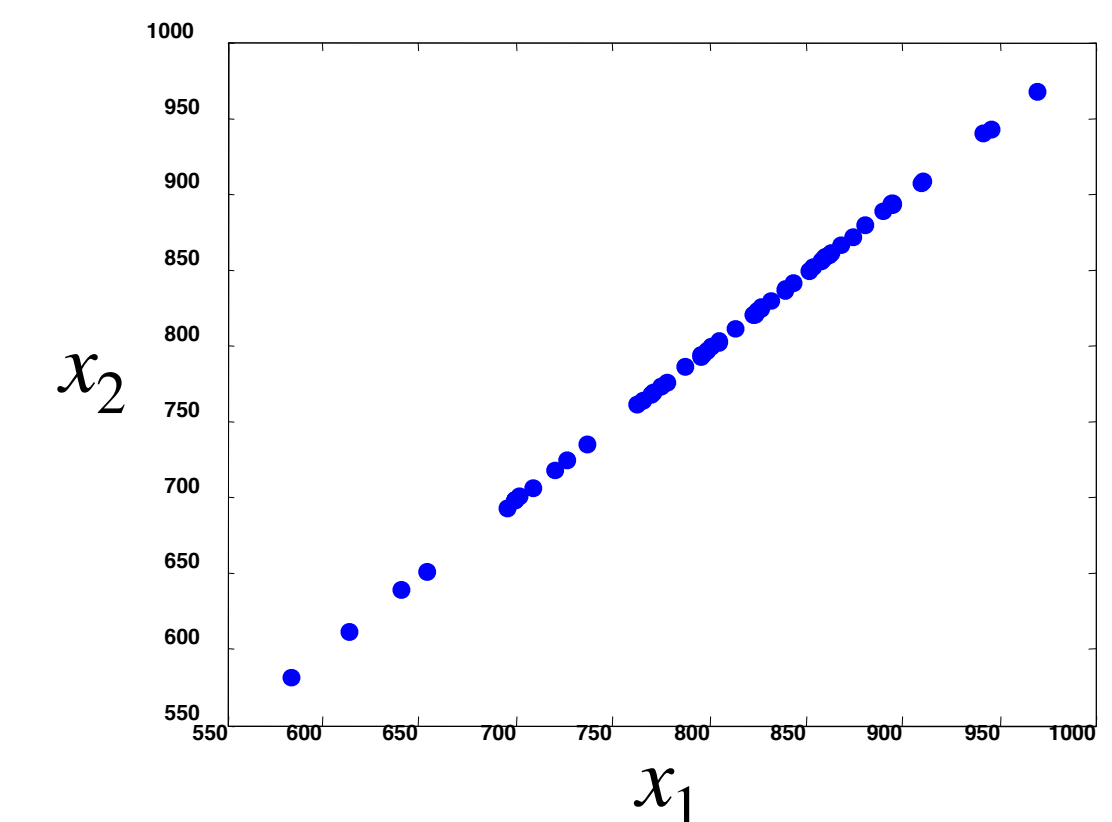
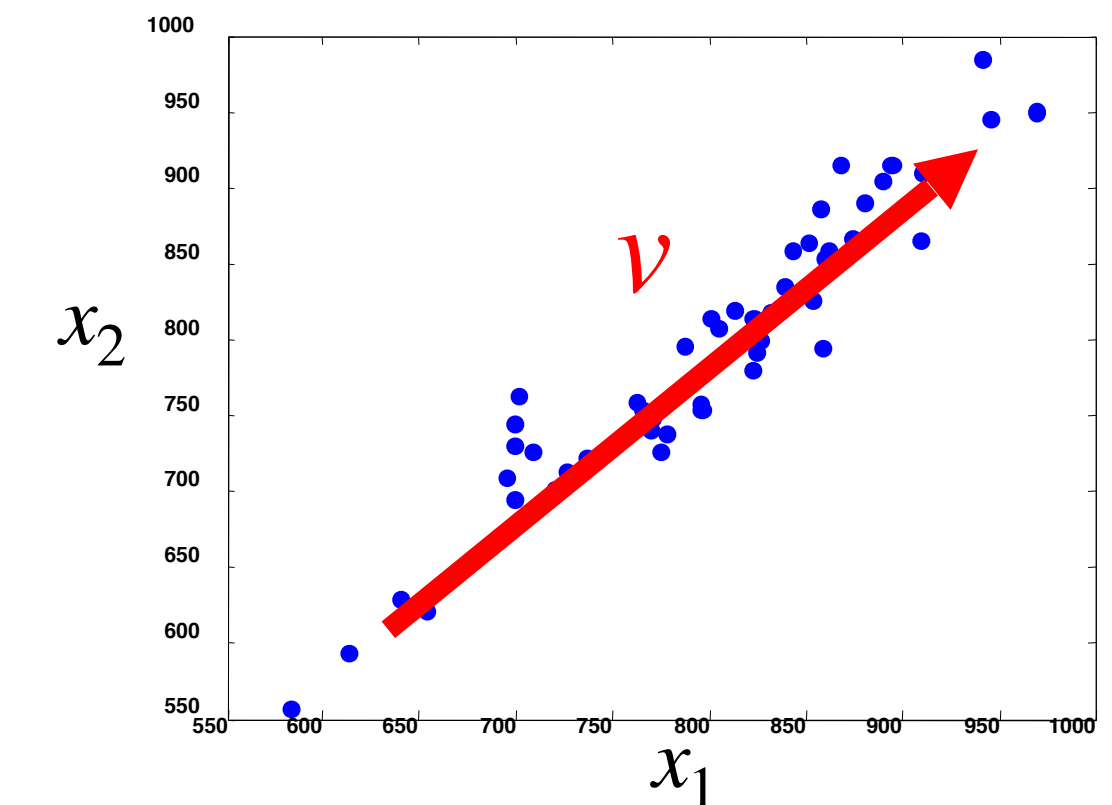
Linear features

- Example: **summarize** two real features $x = [x_1, x_2]$ \rightarrow one real feature z
 - If z **preserves** much information about x , should be able to find $x \approx f(z)$

- **Linear embedding:**

- $x \approx z\nu$
- $z\nu$ should be the **closest** point to x along ν

- $\implies z = x \cdot \nu$



Principal Component Analysis (PCA)

- How to find a good v ?

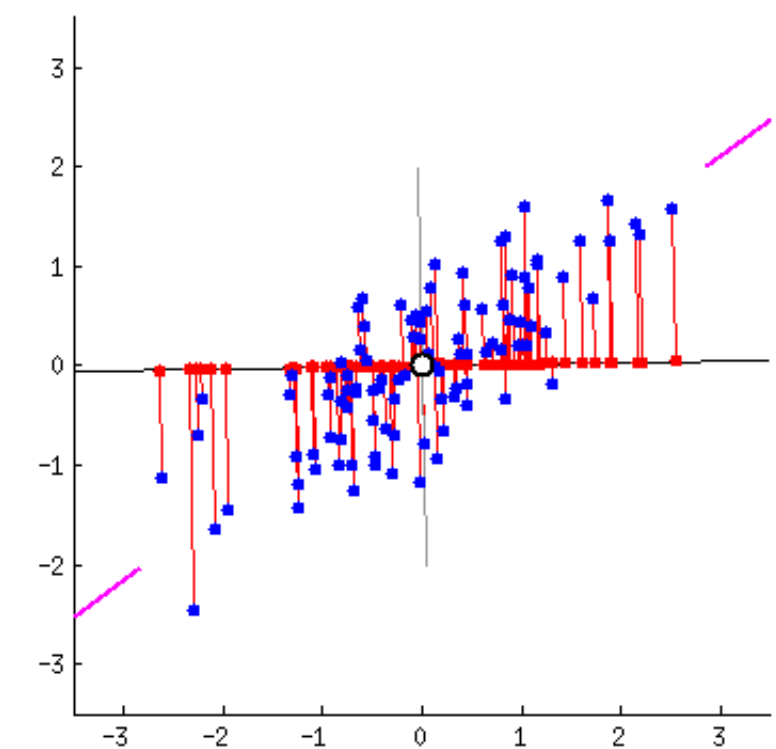
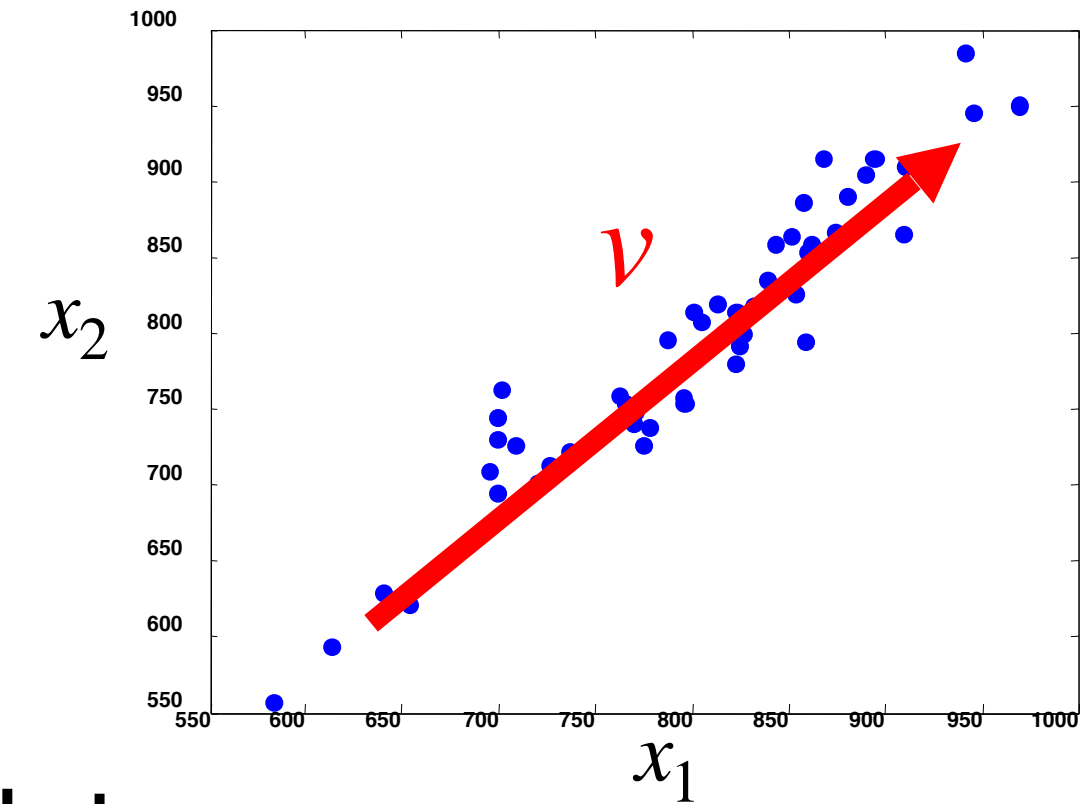
- ▶ Assume X has mean 0; otherwise, subtract the mean $\tilde{X} = X - \mu$
- ▶ Idea: find the direction v of maximum “spread” (variance) of the data

- ▶ Project \tilde{X} on v : $z = \tilde{X}v$

$$\max_{v: \|v\|=1} \sum_i (z_i)^2 = z^T z = v^T \tilde{X}^T \tilde{X} v \implies v \text{ is eigenvector of } \tilde{X}^T \tilde{X} \text{ of largest eigenvalue}$$

empirical covariance

- ▶ = minimum MSE of the residual $\tilde{X} - z v^T = \tilde{X} - \tilde{X} v v^T$



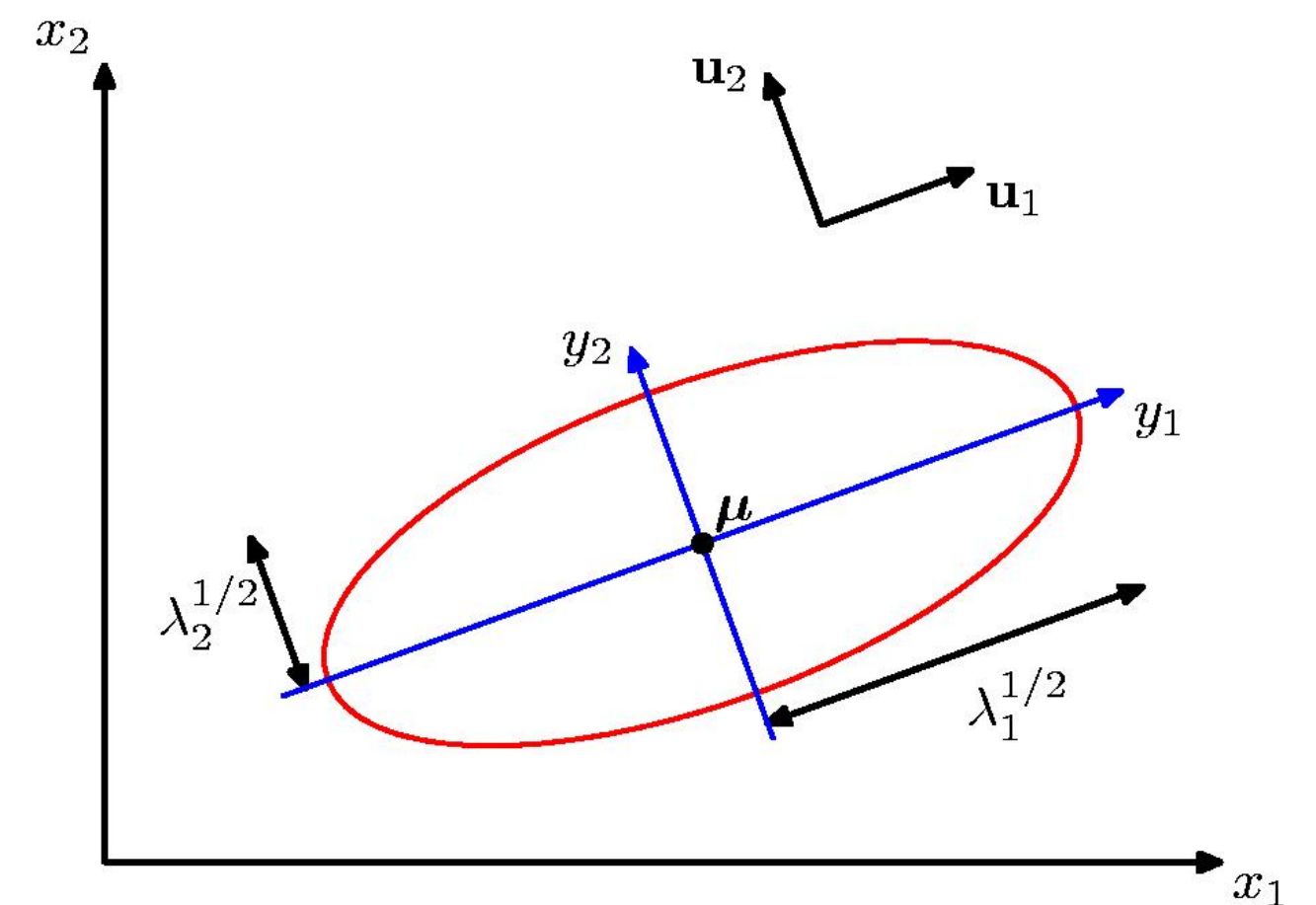
Source

Geometry of a Gaussian

- Data covariance: $\Sigma = \frac{1}{m} \tilde{X}^\top \tilde{X}$ $\tilde{X} = X - \mu$
- Gaussian fit: $p(x) \sim \mathcal{N}(\mu, \Sigma)$
- Value contour for $p(x)$: $\Delta^2 = (x - \mu)^\top \Sigma^{-1} (x - \mu) = \text{const}$
- It's always possible to write Σ in terms of its eigenvectors U , eigenvalues λ :

$$\Sigma = U \Lambda U^\top = \sum_{i=1}^n \lambda_i u_i u_i^\top \implies \Sigma^{-1} = \sum_{i=1}^n \frac{1}{\lambda_i} u_i u_i^\top$$

$$\text{In the eigenvector basis: } \Delta^2 = \sum_{i=1}^n \frac{y_i^2}{\lambda_i}, \text{ with } y_i = u_i^\top (x - \mu)$$



PCA representation

- Subtract data mean from data points
- (Optional) Scale each dimension by its variance
 - Don't just focus on large-scale features (e.g., +1 mileage \ll +1yr ownership)
 - Focus on correlation between features
- Compute empirical covariance matrix $\Sigma = \frac{1}{m} \sum_i \tilde{x}_i \tilde{x}_i^\top$
- Take k largest eigenvectors of $\Sigma = U \Lambda U^\top$

Singular Value Decomposition (SVD)

- Alternative method for finding covariance **eigenvectors**
 - Has many other uses
- **Singular Value Decomposition (SVD):** $X = UDV^T$
 - U and V (left- and right **singular vectors**) are orthogonal: $U^T U = I$, $V^T V = I$
 - D (**singular values**) is rectangular-diagonal
 - $\Sigma = X^T X = VD^T U^T U D V^T = V(D^T D) V^T$
- UD matrix gives **coefficients** to reconstruct data: $x_i = U_{i,1} D_{1,1} v_1 + U_{i,2} D_{2,2} v_2 + \dots$
 - We can truncate this after **top k singular values** (square root of eigenvalues)

$$\begin{array}{|c|} \hline X \\ \hline m \times n \\ \hline \end{array} \approx \begin{array}{|c|} \hline U \\ \hline m \times k \\ \hline \end{array} \cdot \begin{array}{|c|} \hline D \\ \hline k \times k \\ \hline \end{array} \cdot \begin{array}{|c|} \hline V^T \\ \hline k \times n \\ \hline \end{array}$$