# Data Science Survival Skills

Exercise 4 – Data exploration and visualization

# Today's workflow

Data exploration → Create plots → Create figure & edit
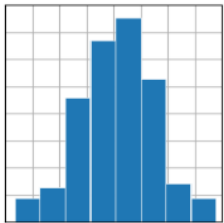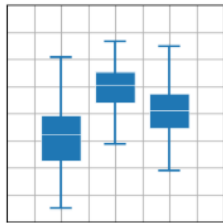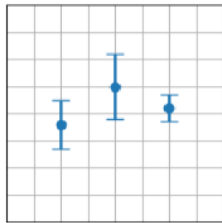
pandas

matplotlib
+
seaborn
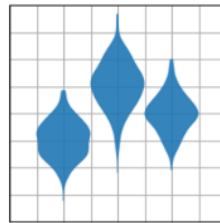
INKSCAPE

# Plotting distributions



hist(x)

boxplot(X)

errorbar(x, y, yerr, xerr)

violinplot(D)

**matplotlib**

histplot(df, x)

boxplot(df, x)

pointplot(df, x)

violinplot(df, x)

**seaborn**

DataFrame

```python
import seaborn as sns
sns.histplot(df, x="Age")
```
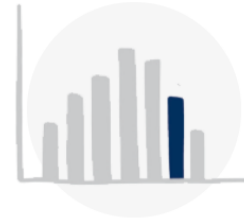
# Visualization



understand the
**context**

choose an
**effective visual**
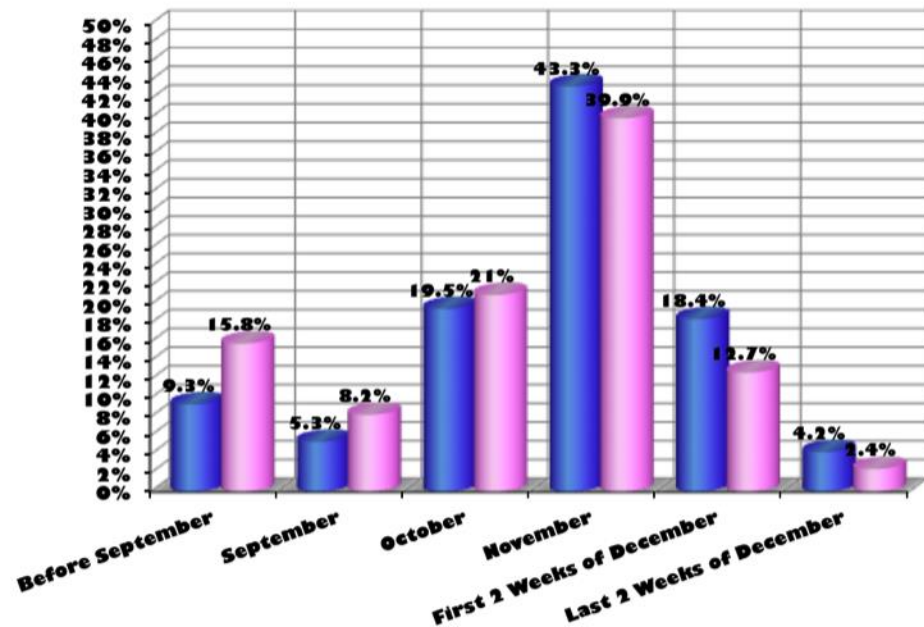
eliminate
**clutter**

focus
**attention**

tell a
**story**

# Removing clutter and focusing attention
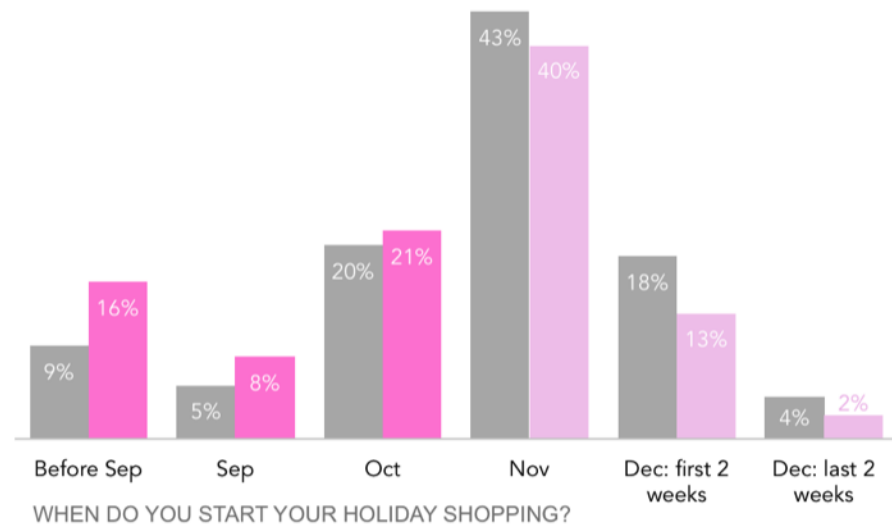


**Shoppers Begins Shopping for Holidays**

Men   Women

| | Men | Women |
|---|---|---|
| Before September | 9.3% | 15.8% |
| September | 5.3% | 8.2% |
| October | 19.5% | 21% |
| November | 43.3% | 39.9% |
| First 2 Weeks of December | 18.4% | 12.7% |
| Last 2 Weeks of December | 4.2% | 2.4% |

**More women start their holiday shopping early**

Men   Women

% OF TOTAL

| | Men | Women |
|---|---|---|
| Before Sep | 9% | 16% |
| Sep | 5% | 8% |
| Oct | 20% | 21% |
| Nov | 43% | 40% |
| Dec: first 2 weeks | 18% | 13% |
| Dec: last 2 weeks | 4% | 2% |

WHEN DO YOU START YOUR HOLIDAY SHOPPING?

# Removing clutter and focusing attention



BEFORE

AFTER

# Some coding

# Raster graphics vs. vector graphics

|  | Raster | Vector |
|---|---|---|
| File extensions | **.png** .jpg .gif | **.svg** .ai .emf |
| Built from … | Pixels | Mathematical equations, lines, and curves |
| Usage | Photos, presentations, web, … | Figures for scientific papers, illustrations, logos, … |
| Pros and cons | - Lose quality when resized<br>+ compatibility | + Don't lose quality when resized (very scalable) |

# Example: Scientific Figure



**Figure 1.** A single latent space channel is sufficient for glottis segmentation. (**A**) Glottis segmentation of endoscopic images using deep neural networks (DNNs) with latent space $\Psi$. (**B**) Convergence of segmentation DNNs across different latent space channels with enabled skip connections. Gradient from black to red indicates fewer channels. The gray line indicates maximum IoU score. (**C**) Convergence of segmentation DNNs across different latent space channels with disabled skip connections. Gradient from black to magenta indicates fewer channels. The gray line indicates maximum IoU score from panel **B**. (**D**) Performance of best performing segmentation DNNs on validation set (solid lines) and evaluated on test set (dashed lines) with enabled (with, red) and disabled (without, magenta) skip connections across latent space ($\Psi$) channels measured by mean intersection over union (IoU). The asterisk indicates the architecture used in the subsequent experiments.

Reference: Kist, Andreas M., et al. "A single latent channel is sufficient for biomedical glottis segmentation." *Scientific Reports* 12.1 (2022): 14292.

# Editing plots in Inkscape

Download: https://inkscape.org/

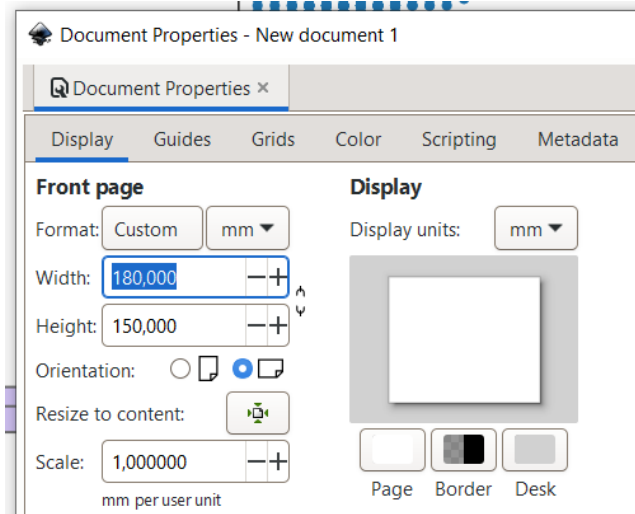Tutorials: https://inkscape-manuals.readthedocs.io/

**Make font editable:**

```python
plt.rcParams['svg.fonttype'] = 'none'
```

# Inkscape



## Set figure size:
File → Document properties



## Snap to grid points

View → ☑ Page Grid

View → Show/Hide → Snap controls bar
→ Activate snapping

## Aligning elements

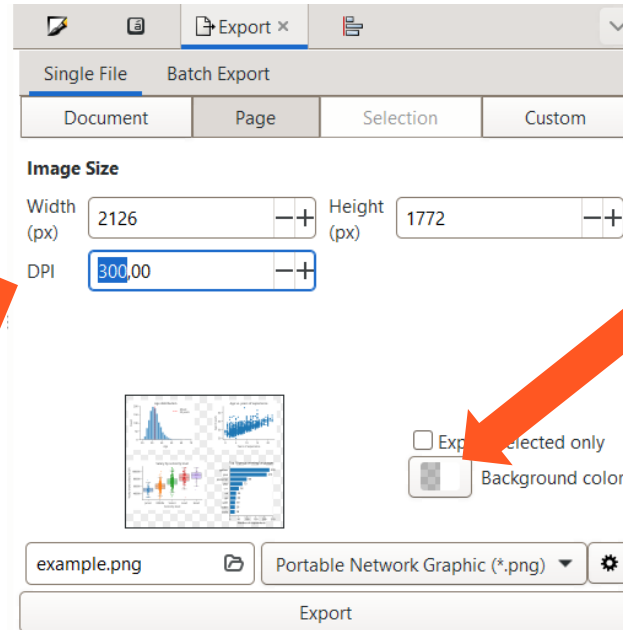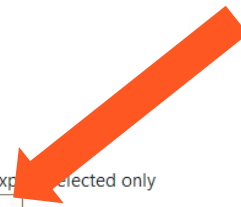Shift + Ctrl + A

# Inkscape

**Save as SVG (vector graphic)**
File → Save as …

**Export as PNG (raster graphic)**
File → Export …



Resolution
(dots per inch)

Set background
to white or
transparent

# More coding

# Questions?

# And we are done!

# Thank you