1. Batch normalization is a technique used in artificial neural networks to improve the training process and performance of deep learning models. It involves normalizing the inputs of each layer in a mini-batch of data, typically right before applying the activation function. This normalization helps stabilize and speed up the training process by reducing internal covariate shift, which refers to the change in the distribution of network activations due to parameter updates during training.

2. There are several benefits of using batch normalization during training:
   - **Faster Convergence**: Batch normalization helps accelerate the training process by reducing the internal covariate shift. This allows the network to converge faster, which means it requires fewer training iterations to reach optimal performance.
   - **Improved Gradient Flow**: By normalizing the inputs to each layer, batch normalization helps to ensure that gradients propagate smoothly through the network during backpropagation. This mitigates the vanishing or exploding gradient problem and enables more stable training.
   - **Regularization Effect**: Batch normalization acts as a form of regularization by adding noise to the inputs of each layer. This noise can help prevent overfitting and improve the generalization ability of the model.
   - **Reduced Sensitivity to Initialization**: Batch normalization makes deep neural networks less sensitive to the choice of initialization parameters, which simplifies the process of training and tuning the model.

3. The working principle of batch normalization involves two main steps: normalization and learnable parameters.
   - **Normalization Step**: In the normalization step, the input to each layer (activation values) is normalized to have a zero mean and unit variance across the mini-batch. Mathematically, this can be expressed as: $\hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}}$ where $x$ is the input to the layer, $\mu$ is the mean of the mini-batch, $\sigma^2$ is the variance of the mini-batch, and $\epsilon$ is a small constant added for numerical stability.

- **Learnable Parameters**: In addition to normalizing the inputs, batch normalization introduces learnable parameters: scale ($\gamma$) and shift ($\beta$) parameters. These parameters allow the network to learn the optimal scaling and shifting of the normalized inputs, enabling it to adapt to the specific requirements of each layer. The normalized input $\hat{x}$ is then transformed as follows: $y = \gamma \hat{x} + \beta$ where $y$ is the output of the batch normalization layer. During training, the scale and shift parameters ($\gamma$ and $\beta$) are learned through backpropagation along with the other network parameters. During inference, batch normalization uses the aggregated statistics (mean and variance) computed during training to normalize the inputs.